

# Verifying a Synthesized Implementation of IEEE-754 Floating-Point Exponential Function using HOL

BEHZAD AKBARPOUR<sup>1,\*</sup>, AMR T. ABDEL-HAMID<sup>2</sup>, SOFIÈNE TAHAR<sup>3</sup>  
AND JOHN HARRISON<sup>4</sup>

<sup>1</sup>*University of Cambridge, Computer Laboratory, William Gates Building, 15 JJ Thomson Avenue, Cambridge CB3 0FD, UK*

<sup>2</sup>*German University of Cairo, Faculty of Information and Electronics Technology, Tagamoa El-Khamis, Cairo, Egypt*

<sup>3</sup>*Concordia University, Department of Electrical and Computer Engineering, 1455 De Maisonneuve Blvd. West, Montreal, Quebec, Canada, H3G 1M8*

<sup>4</sup>*Intel Corporation, JF1-13, 2111 NE 25th Avenue, Hillsboro, Oregon, OR 97124, USA*

*\*Corresponding author: behzad.akbarpour@cl.cam.ac.uk*

**Deep datapath and algorithm complexity have made the verification of floating-point units a very hard task. Most simulation and reachability analysis verification tools fail to verify a circuit with a deep datapath like most industrial floating-point units. Theorem proving, however, offers a better solution to handle such verification. In this paper, we have hierarchically formalized and verified a hardware implementation of the IEEE-754 table-driven floating-point exponential function algorithm using the higher-order logic (HOL) theorem prover. The high ability of abstraction in the HOL verification system allows its use for the verification task over the whole design path of the circuit, starting from gate-level implementation of the circuit up to a high-level mathematical specification.**

*Keywords: floating-point arithmetic; formal hardware verification; higher-order logic; theorem proving*

*Received 8 January 2008; revised 16 March 2009*

*Handling editor: Iain Stewart*

## 1. INTRODUCTION

The verification of floating-point circuits has always been an important part of processor verification. The importance of arithmetic circuit verification was illustrated by the famous floating-point division bug in Intel's Pentium® processor [1]. Floating-point algorithms are usually very complicated. They are composed of many modules where the smallest flaw in design or implementation can cause a very hard to discover bug, as happened in the Intel case. Traditional approaches to verifying floating-point circuits are based on simulation. However, these approaches cannot exhaustively cover the input space of the circuits. In contrast, using formal methods [2] for verification of the correctness of hardware, sometimes just called hardware verification, the behavior of the hardware design is described mathematically, and a formal proof is used to verify that it meets rigorous specifications of intended behavior.

However, formal verification is not the golden rule in circuit testing because of some limitations. A correctness proof cannot guarantee that the real device will never malfunction; the design model of the device may be proved correct, but the hardware actually built can still behave in a way unintended by the designer (this is the case for simulation too). Wrong specification can play a major role in this, because it has been verified that the system will function as specified, but it has not been verified that it will work correctly. Defects in physical fabrication can cause this problem too. In formal verification, a model of the design is verified, not the real physical implementation. Therefore, a fault in the modeling process can give false negatives (errors in the design which do not exist). Because of these limitations, we can consider simulation and formal verification as complementary techniques, the methods have to play together.

Formal verification can be generally divided into two main categories [3]: reachability analysis and deductive methods. Model checkers and equivalence checkers are examples of the first approach. Many different theorem provers (as HOL [4]) have been used for deductive verification. To verify floating-point arithmetic circuits, model checkers would encounter some difficulties as noted in [5]. First, the specification languages are not powerful enough to express arithmetic properties; for arithmetic circuits, the specifications must be expressed as Boolean functions, which is not suitable for complex circuits. Second, these model checkers cannot represent arithmetic circuits efficiently in their models. It is hence no surprise that most related work in the area of formal specification and verification of floating-point arithmetic circuits were done using theorem proving.

Formal verification methods [3] have sometimes been accused of a lack of ability to get into a whole industrial product design cycle. Working on the same design path of most electronic products, we discuss in this paper the formalization and verification of the IEEE-754 [6] table-driven exponential function in all abstraction levels of the design flow. The IEEE-754 exponential function was first specified formally by Harrison [7]. This behavioral specification was written in a high-level *while language*, and was intended mainly to be verified against a more abstract mathematical description of the exponential function [8]. Starting from this behavioral specification Bui *et al.* [9] developed an register-transfer level (RTL) implementation of the design using VHDL and Verilog. In a previous paper [10], Abdel-Hamid *et al.* have introduced design changes to the code produced by Bui *et al.*, to be able to verify this code. They have developed a modular specification and verified the same module, yet this modular specification failed to connect easily to the higher-level algorithmic specification developed by Harrison. The goal of this work is to use formal methods in modeling and verification of the synthesized table-driven exponential function gate-level implementation against the higher-level algorithmic model previously developed by Harrison. In this exercise, we extend Harrison's verification of the exponential function [7] performed as an error analysis between *real* and *algorithmic* levels, first to *RTL* and then to *gate* level, therefore closing the gap between these levels. In contrast to [10] that reconstructed the RTL implementation and established a modular proof between the RTL and behavioral level, we propose a direct verification methodology without any changes to the lower level designs. We will explain in detail how the verification of the synthesized gate level and RTL designs is linked to the algorithmic level for each and every module in the system.

In this work, we use the higher-order logic (HOL) theorem proving system [4] for specifying and verifying the floating-point design at hand. The HOL theorem prover is an interactive proof assistant for HOL developed at Cambridge university by Gordon *et al.* [4]. It was explicitly designed for the formal verification of hardware, though it has also been applied

to other areas including software verification and formalization of pure mathematics. To the best of our knowledge, this is the first attempt to close the verification gap between abstract mathematical specification and a synthesized gate-level implementation using one single formalism and tool, namely HOL.

The organization of the paper is as follows: Section 2 gives a review on work related to the formalization and verification of floating-point algorithms and designs, some of which directly influenced our work. Section 3 describes the table-driven exponential function algorithm, the formal specification and implementation of which are discussed throughout this paper. Section 4 introduces our modeling and verification methodology and shows the main goal we are trying to reach. Section 5 shows the formalized specification of the exponential function in HOL. It also describes the VHDL implementation of the algorithm and introduces its HOL formalization. Section 6 describes the formal verification of the exponential function. We first describe the verification of the exponential function in the transition from the algorithmic level to RTL, using one of the building blocks, namely the floating-point multiplication. The details of the algorithmic to RTL verification of other blocks such as floating-point addition or rounding are given in Appendix A. We then describe the verification of the exponential function in the transition from RTL to gate level, using one of the primitive building blocks, namely the *n-bit Multiplier*. The details of the RTL to gate-level verification of other blocks such as *n-bit Adder* and *n-bit Shifter* are given in Appendix B. Finally, conclusions are drawn in Section 7.

## 2. RELATED WORK

There exist several related work in the open literature on the formalization and verification of IEEE standard-based floating-point arithmetic. For instance, Barrett [11] specified parts of the IEEE-754 standard in Z and Miner [12] formalized the IEEE-854 [13] floating-point standard in PVS. The latter was one of the earliest on the formalization of floating-point standards using theorem proving. This formal specification was then used by Miner and Leathrum [14] to verify in PVS a general class of IEEE-compliant subtractive division algorithms. Carreno [15] formalized the same IEEE-854 standard in HOL. He interpreted the lexical descriptions of the standard into mathematical conditional descriptions and organized them in tables, which were then formalized in HOL.

The most related work among these efforts, however, is the one of Harrison [16] who constructed the real numbers in HOL. He then developed in HOL a generic floating-point library [17] to define the most fundamental terms of the IEEE-754 standard and to prove the corresponding correctness analysis lemmas. He used this library to formalize and verify floating-point algorithms of complex arithmetic operations such as the square root, the exponential function [7] and the transcendental functions [18] against their abstract mathematical counterparts. He also used the floating-point library for the

verification of the class of division algorithms used in the Intel IA-64 architecture [19].

In [20], Moore *et al.* verified the AMD-K5 floating-point division algorithm using the ACL2 theorem prover. Also, Russinoff [21] has developed a floating-point library for the ACL2 prover and applied it successfully to verify the floating-point multiplication, division and square root algorithms of the AMD-K5 and AMD Athlon processors.

In most of the work described above, the scope of the researchers was concentrated in two main fields: first, the formalization of the IEEE floating-point standards and the verification of their relations to the unbounded real numbers as in [12,15,16]; second, the behavioral modeling of floating-point algorithms and verifying their correctness against their main mathematical models as in [7,18].

In [22], Leaser and O'Leary verified a radix-2 square root algorithm and its hardware implementation, used in many processors such as HP PA7200 and Intel Pentium® [21]. They used theorem proving to bridge the abstraction gap between the algorithm and the implementation. The Nuprl proof development system was used for proof automation. This work discusses the proof of the above algorithm starting from RTL and progressing down to gate-level implementation.

Another approach for verification is combining a theorem prover with a model checker or a simulation tool, where the theorem prover handles the high-level proofs, while the low-level properties are handled by the model checker or simulation. For instance, Aagaard and Seger [23] used the Voss hardware verification system to verify the IEEE compliance of a floating-point multiplier. O'Leary *et al.* [24] reported on the specification and verification of the Intel Pentium® Pro processor's floating-point execution unit at the gate level using a combination of model checking and theorem proving. Chen and Bryant [25] used word-level SMV to verify a floating-point adder. Cornea-Hasegan [26] used iterative approaches and mathematical proofs to verify the correctness of the IEEE floating-point square root, divide and remainder algorithms. Compared with theorem proving, this approach is much more automatic, but still requires user guidance.

More recently, Daumas *et al.* [27] have presented a generic library for reasoning about floating-point numbers within the Coq system. This library was then used in the verification of IEEE-compliant floating-point arithmetic algorithms [28] and hardware units [29]. Berg and Jacobi [30] have formally verified a theory of IEEE rounding presented in [31] using the theorem prover PVS. This theory was then used to prove the correctness of a fully IEEE-compliant floating-point unit used in the VAMP processor [32]. Sawada and Gamboa [33] formally verified the correctness of a floating-point square root algorithm used in the IBM Power4™ processor. The verification was carried out with the ACL2(r) theorem prover. Kaivola and coworkers [34–36] presented the formal verification of the floating-point multiplication, division and square root units of the Intel IA-32 Pentium® 4 microprocessor. The verification was carried out

using the Forte verification framework. Both the IBM and Intel floating-point verification efforts use symbolic simulation (via ACL2 at IBM and STE (symbolic trajectory evaluation) at Intel) for verification of optimized gate-level designs against clean RTL models. The automation provided by symbolic simulation is a necessity to keep the amount of human effort down to a reasonable level. However, in our case, it is difficult to describe and verify mathematical circuits using automated tools except for a very limited set of the generated subgoals, therefore we opted for using HOL to solve all different goals interactively. Nevertheless, the produced proof is highly modular and this would allow people to use it as a general framework and change the verification method safely for some of such subgoals. On top of that, we want to link the correctness proof of the RTL to gate-level transition, to the correctness proof of the algorithmic to RTL transition, and also to the error analysis between real and algorithmic levels, and prove a single theorem that connects the floating-point exponential function at the gate level to its abstract mathematical counterpart.

In summary, most openly available related work, except for [22], discuss details of the verification of a hardware implementation, usually at RTL, against predefined properties for the IEEE floating-point standard. This may cover compatibility of the floating-point implementations under investigation to the IEEE standard, but it would not cover the correctness of the implementation against the main circuit behavioral specification. Also, it can be noticed that most of these works are either concerned with the verification of the abstract mathematical description of an IEEE floating-point standard, or is only concerned with the RTL verification against a higher behavioral specification. In this work, we will discuss the formalization and verification of the IEEE-754 table-driven exponential function in all abstraction levels of the design flow.

### 3. THE IEEE-754 EXPONENTIAL FUNCTION ALGORITHM

In this section, we give an introduction to the IEEE-754 exponential function algorithm formal specification and design of which are discussed in the rest of the paper.

Using an approximate polynomial expansion, Tang [8] has developed an algorithm for computing the floating-point exponential function using what he calls a *table-driven* approach. In this approach, given an input argument  $x$ , exceptional cases such as NaN (not-a-number), infinities (or simply very large arguments) and zeros are dealt with first. For example,  $\exp(-\infty) = +0$ . Furthermore, if the argument  $x$  is small enough for this to be a satisfactory approximation, the exponential function is calculated simply as  $1+x$ . The main part of the algorithm covers the remaining cases. Mathematically, the procedure is simple. First we obtain a reduced argument  $r$  such that for some integer  $n$ :

$$x = n \frac{\ln(2)}{32} + r$$

and  $-\ln(2)/64 \leq r \leq \ln(2)/64$ . This  $n$  is found by rounding  $x(32/\ln(2))$  to the nearest integer. Now we decompose  $n$  into its quotient and remainder when divided by 32, i.e.  $n = 32m + j$  with  $0 \leq j \leq 31$ . Hence

$$e^x = e^{(32m+j)(\ln(2)/32)+r} = e^{\ln(2)m} e^{\ln(2)j/32} e^r = 2^m 2^{j/32} e^r$$

Values of  $2^{j/32}$  for  $0 \leq j \leq 31$  are pre-stored constants, and multiplication by  $2^m$  is fast. Hence we just need to calculate  $e^r$  for  $r \in [-\ln(2)/64, \ln(2)/64]$ . This is done by a lower-order polynomial approximation  $p(r) \approx e^r - 1$ , where:

$$p(r) = r + \frac{8388676}{2^{24}} r^2 + \frac{11184876}{2^{26}} r^3$$

The actual reconstruction of  $e^x$ , for reasons of accuracy, is done by:

$$e^x = 2^m (2^{j/32} + 2^{j/32} p(r))$$

In fact, in order to achieve good accuracy, the above mathematical description is complicated slightly. The value  $r$  is broken down into  $r_1 + r_2$ , where  $r_2 \ll r_1$ . Similarly the pre-stored constants  $2^{j/32}$  are all stored as two separate arrays  $S_{\text{lead}}$  and  $S_{\text{trail}}$  with  $2^{j/32} \approx S_{\text{lead}}(j) + S_{\text{trail}}(j)$  and  $S_{\text{trail}}(j) \ll S_{\text{lead}}(j)$ . This would avoid rounding errors as well as take care of the ordering of operations, hence making the actual code look a bit more complicated than the above mathematical description.

#### 4. MODELING AND VERIFICATION METHODOLOGY

The verification process for the table-driven floating-point exponential function will be performed on many levels. Harrison [7] formalized and verified using the HOL Light theorem prover that a behavioral specification of the IEEE-754 table-driven floating-point exponential function implies its abstract mathematical counterpart. He also performed an error analysis between these two levels. For this, he first developed theories in HOL on construction of real numbers [16], and formalization of IEEE-754 standard-based floating-point arithmetic [7, 17]. Then he used valuation functions to find the real value of the floating-point exponential function output, and defined the error as the difference between this value and the corresponding output of the ideal real exponential function. Then he established fundamental lemmas on error analysis of floating-point rounding and arithmetic operations against their abstract mathematical counterparts. Finally based on these lemmas, he proved that the floating-point exponential function algorithm has the correct overflow behavior and, in the absence of overflow, the error in the result is less than 0.54 units in the last place compared with the exact mathematical exponential function. He confirmed and strengthened the main results of the previously published error analysis in [8], though he uncovered a minor error in the hand proof and located a few subtle corners in the proof that a less careful worker might

easily have overlooked. The error in postulated theorems was related to the forgetting of special or degenerate cases in IEEE floating-point such as NaNs and negative zeros.

After handling the transition from real to floating-point levels, we move to the RTL design. At this point, we use the standard HOL predicate approach to model the floating-point exponential function at the RTL, as developed by Bui *et al.* [9] using VHDL and Verilog, within the HOL environment. The last step is to verify this level using a classical hierarchical proof approach in HOL [37]. In this way, we hierarchically prove that the floating-point exponential function RTL implementation implies the high-level algorithmic specification that has already been related to the ideal real specification through the error analysis. The verification can be extended in HOL, following a similar approach, down to gate-level netlist implementation, machine synthesized using the Synopsys tool.

The overall modeling and verification process is described in Fig. 1, where the white boxes are the material provided by [7–9], while the shaded ones represent those developed in this work.

Let  $X$  be the input variable and  $E$  the corresponding output of the floating-point exponential function at the gate level; then our final goal is:

$$\begin{aligned} \vdash_{thm} \forall X E. FP\_EXP\_GATE(X, E) \implies \\ valof(float(E)) = exp(valof(float(X))) \\ + error(X, E) \wedge abs(error(X, E)) \\ \leq error\_bound(X, E) \end{aligned} \quad (1)$$

Here  $FP\_EXP\_GATE$  is a predicate describing the floating-point exponential function in gate level, and its input and output signals  $X$  and  $E$  are Boolean words. To relate these signals to the corresponding specifications in floating-point and real domains, we make use of the bijection function  $float$ , and the valuation function  $valof$ . Also,  $exp$  is the exponential function in real domain available in HOL transcendental functions theory (*transc*). The theorem states that the real value of the floating-point exponential function in gate level is equal to the real value of the exponential function in real domain plus an error, and also the absolute value of the error is bounded to a certain value that depends on the range of the input and output numbers. This goal cannot be reached directly, due to the very high abstraction gap between the gate and abstract mathematics levels as described above. Therefore, the proof scheme was changed to hierarchically prove that the gate level implies the more abstract RTL. Then this RTL was related, by a formal proof, to the behavioral specification. The latter was proved to imply the high-level real specification plus the error. This can be formalized as follows in HOL:

$$\vdash_{thm} \forall X E. FP\_EXP\_GATE(X, E) \implies FP\_EXP\_RTL(X, E) \quad (2)$$

$$\vdash_{thm} \forall X E. FP\_EXP\_RTL(X, E) \implies FP\_EXP\_ALGORITHM(float(X), float(E)) \quad (3)$$

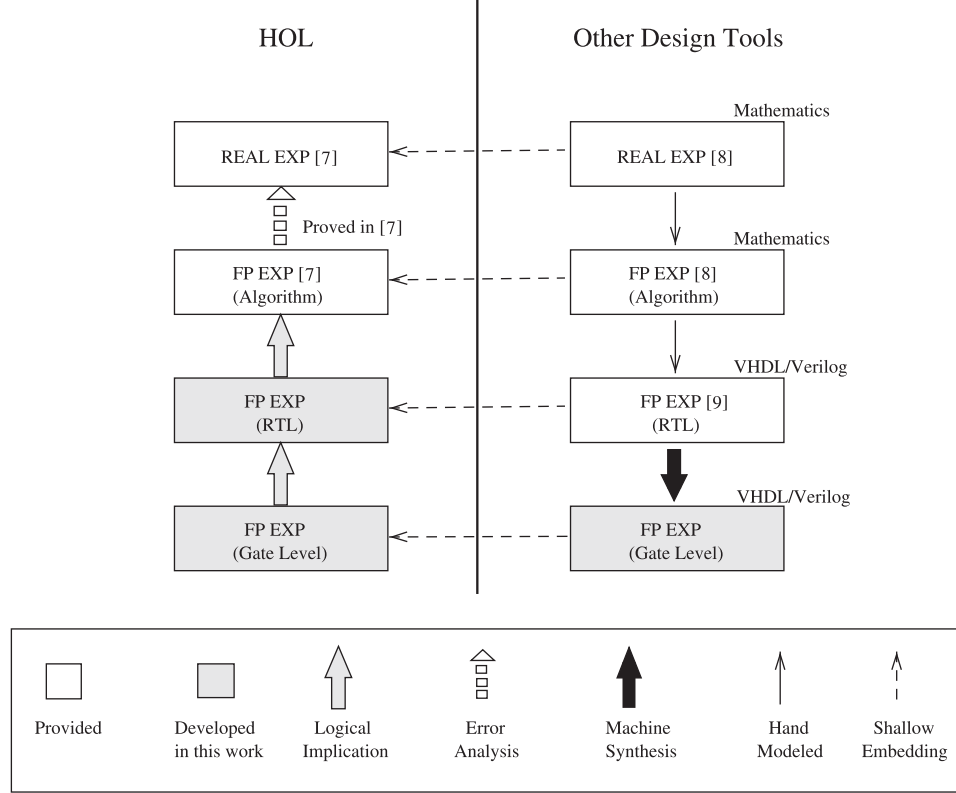


FIGURE 1. Overview of the specification and verification methodology.

$$\begin{aligned}
 & \vdash_{thm} \forall X E. FP\_EXP\_ALGORITHM \\
 & (float(X), float(E)) \implies \\
 & valof(float(E)) = exp(valof(float(X))) \\
 & + error(X, E) \wedge abs(error(X, E)) \\
 & \leq error\_bound(X, E)
 \end{aligned} \tag{4}$$

In these formulas,  $FP\_EXP\_RTL$  and  $FP\_EXP\_ALGORITHM$  are predicates describing the floating-point exponential function in RTL and algorithmic levels, respectively. Note that the inputs and outputs in RTL are still Boolean, however, at the algorithmic level they have floating-point type and we use the data conversion function  $float$  to convert the variables from the Boolean type to IEEE-754 standard-based floating-point type. Also, as can be understood from the theorems, there are no finite precision effects in the transition from gate level to RTL, and also from the RTL to algorithmic level; therefore, the corresponding correctness theorems are described as purely logical implications. However, for the transition from the algorithmic level to the abstract mathematical real number domain, we should consider the effects of finite precision between floating-point numbers and real numbers and conduct an error analysis to bound the corresponding error. Finally using Equations (2–4), we can reach the final goal stated in Equation (1).

Due to the high modularity of the design, the goals of Equations (2) and (3) could be extended to the specification and implementation of sublevel modules, and then the verification continues with these sublevel modules. These proofs were then composed to yield the original goals.

## 5. FORMAL SPECIFICATION AND IMPLEMENTATION OF THE EXPONENTIAL FUNCTION

In this section we describe the formal specification and implementation of the IEEE-754 floating-point exponential function in the HOL theorem prover. The verification details will be discussed in the next section.

### 5.1. Formal specification of the exponential function

The original analysis of the floating-point exponential function in the algorithmic level was performed by Harrison [7] using the HOL Light theorem prover. In this work, we ported the code from HOL Light to HOL4, Kananaskis-4. We modeled the algorithmic specification of the floating-point exponential function as a predicate in HOL as follows:



```

 $\vdash_{def}$  Int_32 = Int(32)
 $\vdash_{def}$  Int_2e9 = Int(2 EXP 9)
 $\vdash_{def}$  Plus_one = float(0,127,0)
 $\vdash_{def}$  THRESHOLD_1 = float(0,134,6056890)
 $\vdash_{def}$  THRESHOLD_2 = float(0,102,0)
 $\vdash_{def}$  Inv_L = float(0,132,3713595)
 $\vdash_{def}$  L1 = float(0,121,3240448)
 $\vdash_{def}$  L2 = float(0,102,4177550)
 $\vdash_{def}$  A1 = float(0,126,68)
 $\vdash_{def}$  A2 = float(0,124,2796268)
 $\vdash_{def}$  FP_EXP_ALGORITHM X E =
   $\exists$  R1 R2 R P Q S E1 N N1 N2 M J S_Lead S_Trail.
    TABLES_OK S_Lead S_Trail  $\wedge$ 
    (if Isnan X then E = X
    else (if X = Plus_infinity then E = Plus_infinity
    else (if X = Minus_infinity then E = Plus_zero
    else (if float_abs X > THRESHOLD_1 then
      (if X > Plus_zero then E = Plus_infinity
      else E = Plus_zero)
    else (if float_abs X < THRESHOLD_2 then E = Plus_one + X
    else
      (N = INTRND (X * Inv_L))  $\wedge$ 
      (N2 = % N Int_32)  $\wedge$ 
      (N1 = N - N2)  $\wedge$ 
      (if Int_abs N  $\geq$  Int_2e9 then
        R1 = X - Tofloat N1 * L1 - Tofloat N2 * L1
      else
        R1 = X - Tofloat N * L1)  $\wedge$ 
      (R2 = Tofloat  $\neg$ N * L2)  $\wedge$ 
      (M = N1 / Int_32)  $\wedge$ 
      (J = N2)  $\wedge$ 
      (R = R1 + R2)  $\wedge$ 
      (Q = R * R * (A1 + R * A2))  $\wedge$ 
      (P = R1 + (R2 + Q))  $\wedge$ 
      (S = S_Lead J + S_Trail J)  $\wedge$ 
      (E1 = S_Lead J + (S_Trail J + S * P))  $\wedge$ 
      E = Scalb (E1,M))))))

```

where the constant TABLES\_OK is used to abbreviate a large set of assumptions about the values of table entries taken from Tang's paper [8]. In addition to IEEE 754 standard single-precision format floating-point numbers, the algorithm uses the formalization of machine integers, which are defined as 2's complement 32-bit integers in HOL.

Based on Tang's algorithm, the above HOL code implements the exponential function in the following four steps:

*Step 1.* Filter out the exceptional cases. When the input argument  $X$  is a NaN, a NaN should be returned. When  $X$  is  $+\infty$ ,  $+\infty$  should be returned without any exception. When  $X$  is  $-\infty$ ,  $+\infty$  should be returned without any exception. When the magnitude of  $X$  is larger than THRESHOLD\_1, a  $+\infty$  with an overflow signal, or a  $+\infty$  with underflow and inexact signals, should be returned. When the magnitude of  $X$  is smaller than THRESHOLD\_2,  $1 + X$  should be returned.

*Step 2.* Reduce the input argument  $X$  to  $[-\frac{\log 2}{64}, \frac{\log 2}{64}]$ . Obtain integers  $M$  and  $J$ , and working-precision floating-point numbers  $R_1$  and  $R_2$  such that (up to roundoff)

$$X = (32M + J) \frac{\log 2}{32} + (R_1 + R_2), \quad |R_1 + R_2| \leq \frac{\log 2}{64}.$$

To perform the argument reduction accurately, do the following:

- Calculate  $N$  as follows:

$$N := \text{INTRND}(X * \text{INV\_L})$$

$$N_2 := N \bmod 32$$

$$N_1 := N - N_2$$

INV\_L is  $\frac{32}{\log 2}$  rounded to working precision. Note that  $N_2 \geq 0$ , regardless of  $N$ 's sign. INTRND rounds a floating-point number to the nearest integer in the manner prescribed by the IEEE standard [6].

- The reduced argument is represented in two working-precision numbers,  $R_1$  and  $R_2$ . We compute them as follows. First, the value of  $\frac{\log 2}{32}$  is represented in two working-precision numbers,  $L_1$  and  $L_2$ , such that the leading part,  $L_1$ , has a few trailing zeros and  $L_1 + L_2$  approximates  $\frac{\log 2}{32}$  to a precision much higher than the working one. If the single-precision exponential is requested and  $|N| \geq 2^9$ , then calculate  $R_1$  by

$$R_1 := (X - N_1 * L_1) - N_2 * L_1.$$

Otherwise, calculate  $R_1$  by

$$R_1 := (X - N * L_1).$$

$R_2$  is obtained by

$$R_2 := -N * L_2.$$

- To complete this step, we decompose  $N$  into  $M$  and  $J$ , thus:

$$\begin{aligned} M &:= \frac{N_1}{32} \\ J &:= N_2. \end{aligned}$$

Step 3. Approximate  $\exp(R_1 + R_2) - 1$  by a polynomial  $p(R_1 + R_2)$ , where

$$p(t) = t + a_1 t^2 + a_2 t^3 + \cdots + a_n t^{n+1}.$$

The polynomial is computed by a standard recurrence:

$$R := R_1 + R_2$$

$$Q := R * R * (A_1 + R * (A_2 + R * (\dots + R * A_n) \dots))$$

$$P := R_1 + (R_2 + Q)$$

The coefficients are obtained from a Remez algorithm implemented by Tang [8]. Our method for bounding the approximation error in this polynomial [7] is *post-hoc*, and

works equally well if the polynomial is derived in other ways, e.g. via Chebyshev expansions [39] or more delicate means [40].

*Step 4.* Reconstruct  $\exp(X)$  via

$$\exp(X) = 2^M (2^{j/32} + 2^{j/32} p(R_1 + R_2)).$$

Each of the values  $2^{j/32}$ ,  $j = 0, 1, \dots, 31$ , is calculated beforehand and represented by two working-precision numbers  $S\_lead(J)$  and  $S\_trail(J)$ . The sum approximates  $2^{j/32}$  to roughly double the working precision. Thus, we may consider  $2^{j/32} = S\_lead(J) + S\_trail(J)$  for all practical purposes. The Reconstruction is as follows:

$$S := S\_lead(J) + S\_trail(J)$$

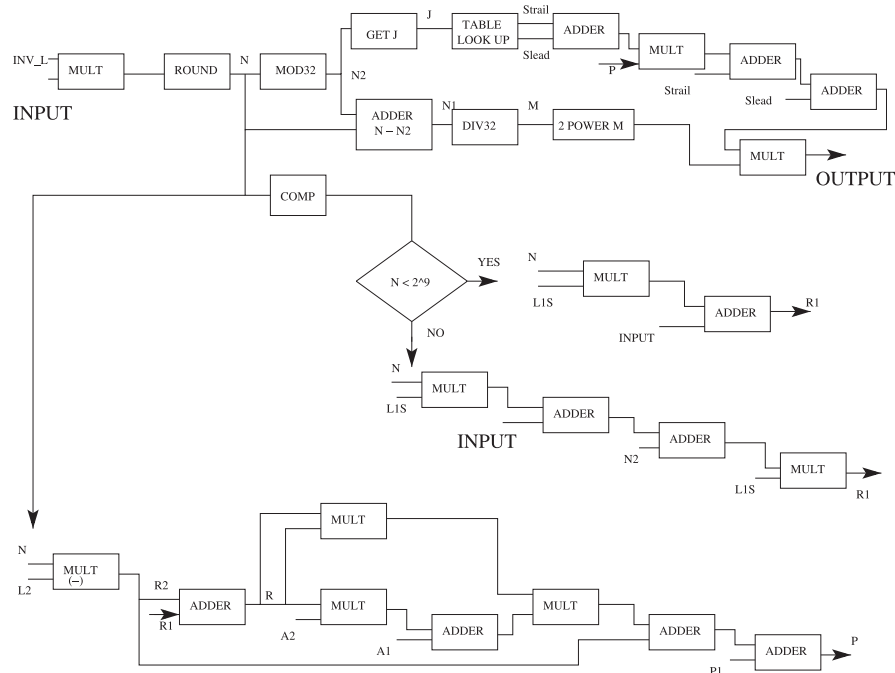
$$\text{exp} := 2^M * (S\_lead(J) + (S\_trail(J) + S * P))$$

## 5.2. Formal implementation of the exponential function

The implementation of the algorithm in RTL was done by Bui *et al.* [9] using two different hardware description languages, namely, Verilog and VHDL.

A block diagram of the whole system is shown in Fig. 2. In this diagram, we use the same labels as in the algorithm specification.

The part constructed using VHDL made use of the sequential mode in contrast to the Verilog implementation that used combinational logic. Both essentially implement the same algorithm outlined in the previous section.



**FIGURE 2.** Floating-point exponential function main block diagram.

The VHDL design is composed of numerous procedures that perform IEEE-754 floating-point operations. These operations include the addition, multiplication, division by 32, rounding to the nearest integer, modulo 32, comparison and powers of 2.

---

```

 $\vdash_{def}$  FP_EXP_RTL xs xe xm outs oute outm =
   $\exists$  inv temp temp2 temp3 twoe9 flag slead strail
    n n1 n2 r1 r2 l1 l2 a1 a2 e1 m q s p r j.
  MULT1 invs xs stemp inve xe etemp invm xm mtemp  $\wedge$ 
  ROUND1 stemp ns etemp ne mtemp nm  $\wedge$ 
  MOD32 ns n2s ne n2e nm n2m  $\wedge$ 
  ADDER1 ns ( $\neg$ n2s) n1s ne n2e n1e nm n2m n1m  $\wedge$ 
  COMP F twoe9s ne twoe9e nm twoe9m flag  $\wedge$ 
  (if flag = WORD [F; F; T] then
    MULT1 ns l1s stemp2 ne l1e etemp2 nm l1m mtemp2  $\wedge$ 
    ADDER1 ( $\neg$ stemp2) xs r1s etemp2 xe r1e xm mtemp2 r1m
  else
    MULT1 ns l1s stemp2 ne l1e etemp2 nm l1m mtemp2  $\wedge$ 
    ADDER1 ( $\neg$ stemp2) xs stemp3 etemp2 xe etemp3 mtemp2 xm mtemp3  $\wedge$ 
    MULT1 stemp2 l1s r1s etemp2 l1e r1e mtemp2 l1m r1m  $\wedge$ 
    ADDER1 ( $\neg$ n2s) stemp3 stemp2 n2e etemp3 etemp2 n2m mtemp3 mtemp2)  $\wedge$ 
  MULT1 ( $\neg$ ns) l2s r2s ne l2e r2e nm l2m r2m  $\wedge$ 
  D32 n1s ms n1e me n1m mm  $\wedge$ 
  ADDER1 r1s r2s rs r1e r2e re r1m r2m rm  $\wedge$ 
  MULT1 rs a2s stemp re a2e etemp rm a2m mtemp  $\wedge$ 
  ADDER1 stemp a1s stemp2 etemp a1e etemp2 mtemp a1m mtemp2  $\wedge$ 
  MULT1 rs rs stemp re re etemp rm rm mtemp  $\wedge$ 
  MULT1 stemp stemp2 qs etemp etemp2 qe mtemp mtemp2 qm  $\wedge$ 
  ADDER1 r2s qs stemp r2e qe etemp r2m qm mtemp  $\wedge$ 
  ADDER1 stemp r1s ps etemp r1e pe mtemp r1m pm  $\wedge$ 
  GET_J n2s n2e n2m j  $\wedge$ 
  TABLES_OK j sleads sleadm sleade  $\wedge$ 
  ADDER1 sleads strails ss sleade straile se sleadm strailm sm  $\wedge$ 
  MULT1 ss ps stemp se pe etemp sm pm mtemp  $\wedge$ 
  ADDER1 stemp strails stemp2 etemp straile se sleadm strailm sm  $\wedge$ 
  ADDER1 sleads stemp2 els sleade etemp2 ele sleadm mtemp2 elm  $\wedge$ 
  TWOPOWERM ms me mm stemp etemp mtemp  $\wedge$ 
  MULT1 stemp els outs etemp ele oute mtemp elm outm

```

---

The design is composed of numerous primitive building blocks including the addition (ADDER1), multiplication (MULT1), division by 32 (D32), rounding to nearest integer (ROUND1), modulo 32 (MOD32), comparison (COMP), powers of 2 (TWOPOWERM) and get J (Get\_J), which will be explained in the next section.

## 6. FORMAL VERIFICATION OF THE EXPONENTIAL FUNCTION

In this section we describe the verification of the floating-point exponential function using HOL according to the methodology described in Section 4. We first describe the verification of the exponential function in the transition from the algorithmic level to the RTL, using one of the building blocks, namely the floating-point multiplication. The details of the algorithmic to RTL verification of other blocks such as floating-point addition, division by 32, round to nearest integer, modulo 32, comparison, powers of two and get J blocks are given in Appendix A. We then describe the verification

To ensure that the code is synthesizable, the program was made primitive and the length was much greater than it needed to be.

We modeled this implementation as a predicate in HOL as follows:

of the exponential function in the transition from the RTL to gate level, using one of the primitive building blocks, namely the n-bit Multiplier. The details of the RTL to gate level verification of other blocks such as n-bit Adder, n-bit Subtractor, n-bit Concatenator, n-bit Multiplexer and n-bit Shifter are given in Appendix B.

### 6.1. Verification of RTL to algorithmic level

In this section we describe the algorithmic level to RTL verification of the floating-point exponential function. The whole RTL design is segmented into different blocks and then modeled using HOL. The resulting model is in turn set against the algorithmic specification and the HOL tool is used interactively to prove its correctness.

*The main theorem.* We established the correctness of the RTL implementation of the floating-point exponential function against its algorithmic specification in HOL as the following main theorem:



Theorem 1: FP\_EXP\_RTL\_TO\_ALGORITHM\_THM

$\vdash \text{FP\_EXP\_RTL } xs \ xe \ xm \ outs \ oute \ outm \implies$   
 $\text{FP\_EXP\_ALGORITHM } (\text{float } (\text{BV } xs, \text{BNVAL } xe, \text{BNVAL } xm))$   
 $(\text{float } (\text{BV } outs, \text{BNVAL } oute, \text{BNVAL } outm))$

where `float` is the bijection function that converts a triplet of natural numbers to the floating-point type, and `BV` and `BNVAL` are predefined functions of the HOL *word* library mapping a single bit and a Boolean word into a natural number, respectively.

---

```
e (REPEAT GEN_TAC THEN
  REWRITE_TAC [FP_EXP_ALGORITHM, FP_EXP_RTL]
  REPEAT STRIP_TAC THEN
  .
  .
  ARW_TAC [MULT1_RTL_TO_ALGORITHM_Correct, ADDER1_RTL_TO_ALGORITHM_Correct,
    D32_RTL_TO_ALGORITHM_Correct, ROUND1_RTL_TO_ALGORITHM_Correct,
    MOD32_RTL_TO_ALGORITHM_Correct, COMP_RTL_TO_ALGORITHM_Correct,
    TWOPOWERM_RTL_TO_ALGORITHM_Correct, GET_J_RTL_TO_ALGORITHM_Correct])
```

---

where lemmas such as `MULT1_RTL_TO_ALGORITHM_Correct`, `ADDER1_RTL_TO_ALGORITHM_Correct`, etc. are about the correctness of the sublevel modules, which relate the RTL implementation of each module with the corresponding algorithmic specification.

In the following sections we will describe in detail the verification of one of the primitive building blocks, namely the floating-point multiplication. The rest is given in Appendix A. For all the blocks described, the RTL descriptions, the corresponding HOL models and parts of the proof strategy are provided to explain the verification in its entirety.

*Verification of floating-point multiplication block.* Multiplication is an operation that is quite straightforward. Its algorithm is divided into three main parts corresponding to the three parts of the single-precision format. The first part, the sign, is determined by an exclusive OR function of the two input signs. The exponent of the output, the second part, is calculated by adding the two input exponents. And finally, the significand is determined

As explained before, there is a high level of regularity and modularity in the design of the floating-point exponential function so that primitive blocks such as adders and multipliers are used to build the larger and complicated design. Also, the main verification goal of the whole design can be broken down to the verification proofs of the sublevel modules. These proofs are then composed to yield the original goals. Therefore the main theorem `FP_EXP_RTL_TO_ALGORITHM_THM` was proved in HOL using the following tactic:

by multiplying the two input significands each with a ‘1’ concatenated to it. The result obtained will have about twice as many bits as the significand should normally have and so, the result will be truncated, normalized and the implied ‘1’ will be removed (see Fig. 3 for the block diagram). The normalization process will be fairly simple knowing that the multiplication of two 24 bit numbers with a one at the most significant bit position will yield a result with a one at the most significant bit (bit 47) or at bit 46. Depending on the situation, the result will either be shifted once or twice. At the beginning of the algorithm, there is an IF statement that checks for exceptional cases where there is a zero in at least one of the inputs. It is important to note that this implementation of the floating-point multiplier does not handle subnormal numbers; therefore, it is not a fully fledged floating-point multiplier. It is a perfect block for the proposed exponential function, as the subnormal numbers are not allowed to reach the multiplier block in this design.

In HOL, we modeled this algorithm as follows:

---

```
 $\vdash_{\text{def}}$  MULT1_RTL s1 s2 s3 e1 e2 e3 m1 m2 m3 =
   $\exists$  imp1 imp2 count mbuff1 mbuff2 mbuff3 mbuff4 mbuff5.
  (if (BNVAL e1 = 0)  $\vee$  (BNVAL e2 = 0) then
    (e3 = NBWORD 8 0)  $\wedge$  (m3 = NBWORD 23 0)  $\wedge$  (s3 = F)
  else
    (s3 = s1 xor s2)  $\wedge$  (mbuff3 = BNVAL e1 - 127)  $\wedge$ 
    (mbuff4 = BNVAL e2 - 127)  $\wedge$  (mbuff2 = mbuff3 + mbuff4)  $\wedge$ 
    (imp1 = WCAT (WORD [T], m1))  $\wedge$  (imp2 = WCAT (WORD [T], m2))  $\wedge$ 
    (mbuff1 = NBWORD 48 (BNVAL imp1 * BNVAL imp2))  $\wedge$ 
    (if BIT 47 mbuff1 = T then count = 1 else count = 2)  $\wedge$ 
    (mbuff5 = SND (SHL F mbuff1 F))  $\wedge$ 
    (m3 = WSEG 23 25 mbuff5)  $\wedge$ 
    (e3 = NBWORD 8 (mbuff2 - count + 127)))
```

---

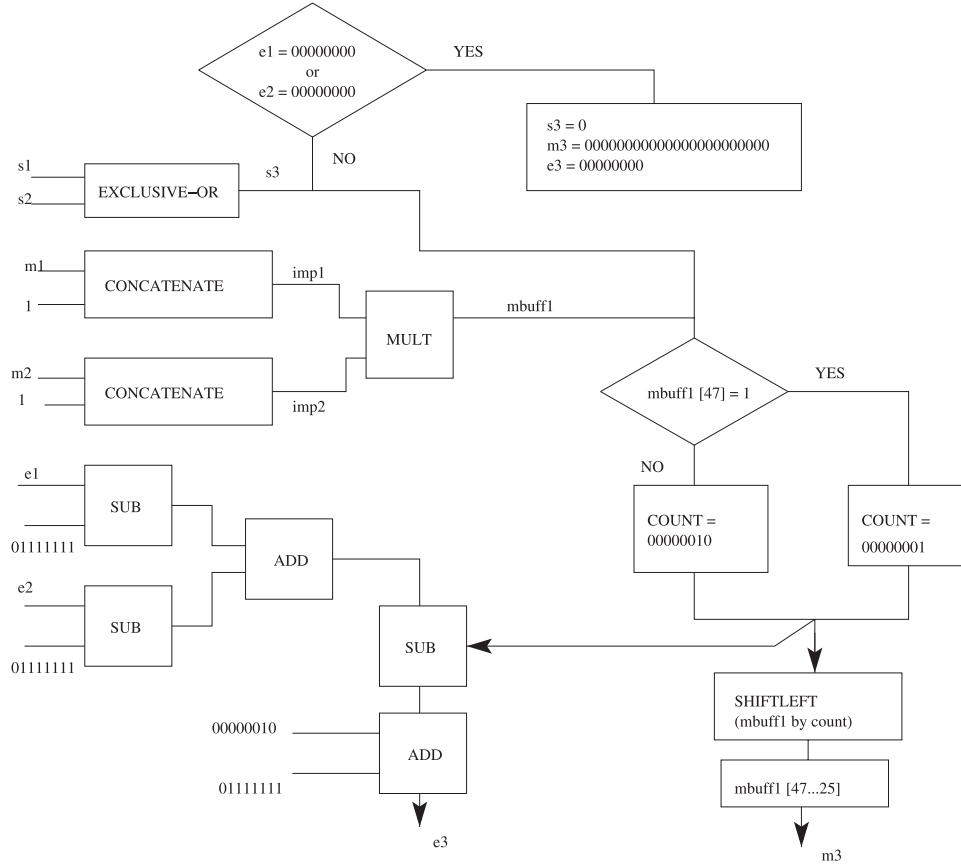


FIGURE 3. Multiplication block diagram.

where BV and BNVAL are predefined functions of the HOL *word* library mapping a single bit and a Boolean word into a natural number, respectively. NEWORD is the reverse function mapping a natural number into a Boolean word with a given word length. BIT, WSEG and WCAT are the basic constants denoting the functions of indexing, segmenting and concatenation of words, respectively, and SHL is the generic shift left operator.

Then we established the correctness of the RTL implementation of the floating-point multiplication function against its algorithmic specification in HOL as the following lemma:

Lemma 1: MULT1\_RTL\_TO\_ALGORITHM\_Correct  
 $\vdash \text{MULT1\_RTL } s1 \ s2 \ s3 \ e1 \ e2 \ e3 \ m1 \ m2 \ m3 \Rightarrow$   
 $(\text{float } (\text{BV } s3, \text{BNVAL } e3, \text{BNVAL } m3)) =$   
 $\text{float } (\text{BV } s1, \text{BNVAL } e1, \text{BNVAL } m1) * \text{float } (\text{BV } s2, \text{BNVAL } e2, \text{BNVAL } m2))$

where float is the bijection function that converts a triplet of natural numbers to the floating-point type. Note that we used the conventional symbols for arithmetic operations on floating-point numbers at the algorithmic level using the operator overloading feature of HOL. The arithmetic operations on floating-point numbers are defined where they first deal with the exceptional cases, either where the arguments involve a NaN

or infinity, or are invalid for other reasons (e.g.  $\infty - \infty$ ) and generate a NaN. Apart from that, they basically just take the real value of the arguments, perform the mathematical operations using the arbitrary precision in real domain and then round the result according to the desired rounding mode. Therefore, our main task in the proof of the above-mentioned theorem was to show that the result of the operation following the RTL algorithm is the best approximation to the real result. These are established in HOL as the following lemmas:

Lemma 2:  
 $\vdash \text{MULT1\_RTL } s1 \ s2 \ s3 \ e1 \ e2 \ e3 \ m1 \ m2 \ m3 \Rightarrow$   
 $((\text{BV } s3, \text{BNVAL } e3, \text{BNVAL } m3) =$   
 $\text{round float\_format To\_nearest}$   
 $(\text{valof float\_format } (\text{BV } s1, \text{BNVAL } e1, \text{BNVAL } m1) * \text{valof float\_format } (\text{BV } s2, \text{BNVAL } e2, \text{BNVAL } m2)))$

Lemma 3:  
 $\vdash \text{MULT1\_RTL } s1 \ s2 \ s3 \ e1 \ e2 \ e3 \ m1 \ m2 \ m3 \Rightarrow$   
 $((\text{BV } s3, \text{BNVAL } e3, \text{BNVAL } m3) =$   
 $\text{closest } (\text{valof float\_format}) (\lambda a. T)$   
 $\{a \mid \text{is\_finite float\_format } a\}$   
 $(\text{valof float\_format } (\text{BV } s1, \text{BNVAL } e1, \text{BNVAL } m1) * \text{valof float\_format } (\text{BV } s2, \text{BNVAL } e2, \text{BNVAL } m2)))$

where round is the floating-point rounding function, float\_format is the floating-point format, To\_nearest is

the rounding to nearest mode, `valof` is the valuation function, `closest` is the function that picks out the best approximation to a real value from a set of floating-point numbers and `is_finite` defines the finiteness criteria for the floating-point numbers. The proof is done by rewriting with the definition of `MULT1_RTL` and a search in the range of all finite floating-point numbers to check if the result of multiplication using this function is the closest value to the real value resulting from multiplication of the real values of two input floating-point numbers.

Following a similar approach, we have verified other building blocks of the floating-point exponential function such as floating-point addition, division by 32, round to nearest integer, modulo 32, comparison, powers of two and get J blocks in the transition from the algorithmic level to RTL using HOL. For more details, please refer to Appendix A.

## 6.2. Verification of gate level to RTL

Following a similar approach to the verification of the RTL to the algorithmic level as described in the previous section, we established the correctness of the gate-level implementation of the floating-point exponential function against its RTL specification in HOL as the following main theorem:

Theorem 2: `FP_EXP_GATE_LEVEL_TO_RTL_THM`  
 $\vdash \text{FP\_EXP\_GATE } xs \ xe \ xm \ outs \ oute \ outm \implies$   
 $\text{FP\_EXP\_RTL } xs \ xe \ xm \ outs \ oute \ outm$

To prove this theorem, we have proved the following lemmas regarding the correctness of each module:

Lemma 12: `MULT1_GATE_TO_RTL_Correct`  
 $\vdash \text{MULT1\_GATE } s1 \ s2 \ s3 \ e1 \ e2 \ e3 \ m1 \ m2 \ m3 \implies$   
 $\text{MULT1\_RTL } s1 \ s2 \ s3 \ e1 \ e2 \ e3 \ m1 \ m2 \ m3$

Lemma 13: `ADDER1_GATE_TO_RTL_Correct`  
 $\vdash \text{ADDER1\_GATE } s1 \ s2 \ s3 \ e1 \ e2 \ e3 \ m1 \ m2 \ m3 \implies$   
 $\text{ADDER1\_RTL } s1 \ s2 \ s3 \ e1 \ e2 \ e3 \ m1 \ m2 \ m3$   
 Lemma 14: `D32_GATE_TO_RTL_Correct`  
 $\vdash \text{D32\_GATE } s1 \ s2 \ e1 \ e2 \ m1 \ m2 \implies$   
 $\text{D32\_RTL } s1 \ s2 \ e1 \ e2 \ m1 \ m2$   
 Lemma 15: `ROUND1_GATE_TO_RTL_Correct`  
 $\vdash \text{ROUND1\_GATE } s1 \ s2 \ e1 \ e2 \ m1 \ m2 \implies$   
 $\text{ROUND1\_RTL } s1 \ s2 \ e1 \ e2 \ m1 \ m2$   
 Lemma 16: `MOD32_GATE_TO_RTL_Correct`  
 $\vdash \text{MOD32\_GATE } s1 \ s2 \ e1 \ e2 \ m1 \ m2 \implies$   
 $\text{MOD32\_RTL } s1 \ s2 \ e1 \ e2 \ m1 \ m2$   
 Lemma 17: `COMP_GATE_TO_RTL_Correct`  
 $\vdash \text{COMP\_GATE } s1 \ s2 \ e1 \ e2 \ m1 \ m2 \implies$   
 $\text{COMP\_RTL } s1 \ s2 \ e1 \ e2 \ m1 \ m2$   
 Lemma 18: `TWOPOWERM_GATE_TO_RTL_Correct`  
 $\vdash \text{TWOPOWERM\_GATE } s1 \ e1 \ m1 \ s3 \ e3 \ m3 \implies$   
 $\text{TWOPOWERM\_RTL } s1 \ e1 \ m1 \ s3 \ e3 \ m3$   
 Lemma 19: `GET_J_GATE_TO_RTL_Correct`  
 $\vdash \text{GET\_J\_GATE } s1 \ e1 \ m1 \ j \implies$   
 $\text{GET\_J\_RTL } s1 \ e1 \ m1 \ j$

The gate-level specification of the modules is very similar to their RTL specification so that they are composed of the same number of sub-modules at the lower level. As can be seen from Figs 3 and A1–A6, there are seven main primitive building block sub-modules in these levels namely n-bit Adder, n-bit Subtractor, n-bit Multiplier, n-bit Comparator, n-bit Concatenator, n-bit Multiplexer, and n-bit Shifter. We use these intermediate sub-modules to cover the gap between the RTL and the gate level. In the following we describe the details of the verification of one such sub-module, n-bit Multiplier. The others are given in Appendix B.

*Verification of n-bit Multiplier.* The n-bit Multiplier in RTL is specified as follows:

---

$\vdash_{\text{def}} \text{CELL\_MUL\_SPEC } a \ b \ c \ p \ co \ po =$   
 $(BV \ po = (\text{if } (BV \ (a \wedge b) + BV \ c + BV \ p < 2) \text{ then}$   
 $(BV \ (a \wedge b) + BV \ c + BV \ p)$   
 $\text{else}$   
 $(BV \ (a \wedge b) + BV \ c + BV \ p) - 2)) \wedge$   
 $(co = \neg (BV \ (a \wedge b) + BV \ c + BV \ p < 2))$

$\vdash_{\text{def}} \text{ShiftLeFT\_Spec } n \ X \ Y = \forall n. ((Y \ 0 = F) \wedge (Y \ (\text{SUC } n) = X \ n))$

$\vdash_{\text{def}} \text{ROW\_MUL\_SPEC } n \ A \ B \ C \ P \ CO \ PO \ Aout =$   
 $\exists c.$   
 $(BV \ (PO \ n) = \text{if } (BV \ ((A \ n) \wedge b) + BV \ (C \ n) + BV \ (P \ n) < 2) \text{ then}$   
 $(BV \ ((A \ n) \wedge b) + BV \ (C \ n) + BV \ (P \ n))$   
 $\text{else}$   
 $(BV \ ((A \ n) \wedge b) + BV \ (C \ n) + BV \ (P \ n)) - 2) \wedge$   
 $((c \ n) = \neg (BV \ ((A \ n) \wedge b) + BV \ (C \ n) + BV \ (P \ n) < 2)) \wedge$   
 $\text{ShiftLeFT\_Spec } n \ A \ Aout \wedge$   
 $\text{ShiftLeFT\_Spec } n \ c \ CO \wedge$   
 $(C \ (\text{SUC } n) = F) \wedge$   
 $(P \ (\text{SUC } n) = F)$

$\vdash_{\text{def}} (\text{ARRAY\_MUL\_spec } 0 \ A \ B \ C \ P \ Co \ Po \ Aout =$   
 $\text{ROW\_MUL\_SPEC } 0 \ A \ (B \ 0) \ C \ P \ Co \ Po \ Aout) \wedge$   
 $(\text{ARRAY\_MUL\_spec } (\text{SUC } n) \ A \ B \ C \ P \ Co \ Po \ Aout = \exists a \ p \ c.$

---

```

ARRAY_MUL_spec n A B C P c p a ^
ROW_MUL_SPEC n a (B (SUC n)) c p Co Po Aout)

```

```

⊢def MUL_SPEC n A B C P MULout =
  ∃ Co Po Aout.
    ARRAY_MUL_spec n A B C P Co Po Aout ^
    nadd_spec ((2 * n) - 1) Co Po F MULout (MULout (2 * n))

```

The  $n$ -bit Multiplier at the gate level is implemented as follows:

```

⊢def CELL_MUL_IMP a b c p co po =
  ∃ s1.
    (and2 a b s1) ^
    (fa_imp s1 c p po co)

⊢def ROW_MUL_IMP n A b C P CO PO Aout =
  ∃ c.
    CELL_MUL_IMP (A n) b (C n) (P n) (c n) (PO n) ^
    ShiftLeFT_Imp n A Aout ^
    ShiftLeFT_Imp n c CO ^
    (C (SUC n) = F) ^
    (P (SUC n) = F)

⊢def ARRAY_MUL_IMP 0 A B C P Co Po Aout =
  ROW_MUL_IMP 0 A (B 0) C P Co Po Aout ^
  (ARRAY_MUL_IMP (SUC n) A B C P Co Po Aout =
   ∃ a p c.
    ARRAY_MUL_IMP n A B C P c p a ^
    ROW_MUL_IMP n a (B (SUC n)) c p Co Po Aout)

⊢def MUL_IMP n A B C P MULout =
  ∃ Co Po Aout.
    ARRAY_MUL_IMP n A B C P Co Po Aout ^
    nadd_imp (2*n-1) Co Po F MULout (MULout (2*n))

```

The correctness of the  $n$ -bit Multiplier block is proved in HOL as in the following theorems:

**Theorem 3: FP\_EXP\_GATE\_LEVEL\_TO\_REAL\_THM**

```

⊢ FP_EXP_GATE xs xe xm outs oute outm ^ Finite (float (BV xs,BNVAL xe,BNVAL xm)) ^
  exp(valof float_format (BV xs,BNVAL xe,BNVAL xm)) < threshold (float_format) ==>
  Isnormal (float (BV outs,BNVAL oute,BNVAL outm)) ^
  abs(valof float_format (BV outs,BNVAL oute,BNVAL outm) -
  exp(valof float_format (BV xs,BNVAL xe,BNVAL xm))) <
  (&54 / &100) * Ulp(float (BV outs,BNVAL oute,BNVAL outm)) ∨
  (Isdenormal(float (BV outs,BNVAL oute,BNVAL outm)) ∨
  Iszero (float (BV outs,BNVAL oute,BNVAL outm))) ^
  abs(valof float_format (BV outs,BNVAL oute,BNVAL outm) -
  exp(valof float_format (BV outs,BNVAL oute,BNVAL outm))) <
  (&77 / &100) * Ulp(float (BV outs,BNVAL oute,BNVAL outm))

```

This main theorem connects the floating-point exponential function at the gate level to its abstract mathematical counterpart. The specification it proves is that the function has the correct overflow behavior and, in the absence of overflow, the error in the result is  $<0.54$  units in the last place (Ulp) (0.77 if the answer is denormalized) compared with the exact mathematical exponential function. One Ulp is defined as the magnitude of the least significant bit of the value concerned.

**Theorem: N\_MUL\_GATE\_LEVEL\_TO\_RTL\_Correct**  
 $\vdash \text{MUL\_IMP } n \ A \ B \ C \ P \ \text{MULout} \implies$   
 $\text{MUL\_SPEC } n \ A \ B \ C \ P \ \text{MULout}$

This goal can be tackled by dividing it into smaller subgoals, where every subgoal represents the verification of one of its sub-modules. This was done by starting with verifying the cell, then the row and then the array multiplier.

Following a similar approach, we have verified other primitive building blocks of the floating-point exponential function such as  $n$ -bit Adder,  $n$ -bit Subtractor,  $n$ -bit Concatenator,  $n$ -bit Multiplexer and  $n$ -bit Shifter in the transition from the RTL to gate level in HOL. For more details, please refer to Appendix B.

### 6.3. Summary

Having proved the Theorems 1 and 2, which state the correctness of the floating-point exponential function in the transition from the gate level to the RTL and algorithmic levels, together with the final correctness theorem proved in [7] about the error analysis of the algorithmic level to real numbers, we can prove the following theorem that bridges the gap between the gate level and ideal real numbers considering the error analysis:

## 7. CONCLUSIONS

Most verification and testing tools will fall short of verifying a circuit with a deep datapath. The IEEE-754 table-driven exponential function with its 32 bit input and 32 bit output implementation would be considered an impossible task for exhaustive simulation. For full coverage with simulation we would have  $2^{32}$  cases, which means that even a 2% or 3%

coverage would take very long simulation time. Model checking techniques will not go a lot further as the deep datapath means a huge state space causing a *state space explosion* [38], making it impossible to verify such a circuit. The properties of the main module and most of its sub-modules cannot be covered easily with, for example, CTL properties [38].

In this paper, we have demonstrated the use of HOL to establish a complete proof between the lower gate level and RTL implementations and the higher-level algorithmic specifications previously developed by Harrison for the IEEE-754 table-driven floating-point exponential function. To establish this proof, we had to formally specify and verify many floating-point smaller modules, such as floating-point addition and floating-point multiplication, as well as many other primitive building blocks. The project was first defined as a two-year master thesis of the second author and then completed by the first author as a half man-year postdoctoral research. The whole code was composed of nearly 5000 lines.

One of the very important advantages of the hierarchical verification lies in the fact that the change of a module or more will not mean the re-proof of the whole system. It only means the re-proof that the new module meets the same specification that the older version did. This may mean a lot for tight time-to-market requirements in a fast-moving technology like electronics. As an example, our proof can always be used with the changing technology as long as we prove that the lower modules, gates for instance, are still satisfying the same properties.

## REFERENCES

- [1] CS-013007 (1994) *Statistical analysis of floating point flaw, description of the flaw*. Intel White Paper, Santa Clara, USA.
- [2] Kropf, T. (2000) *Introduction to Formal Hardware Verification*. Springer, Berlin.
- [3] Kern, C. and Greenstreet, M.R. (1999) Formal verification in hardware design: a survey. *ACM Trans. Des. Autom. Electron. Syst.* **4**, 123–193.
- [4] Gordon, M.J.C. and Melham, T.F. (1993) *Introduction to HOL: A Theorem Proving Environment for Higher-Order Logic*. Cambridge University Press, Cambridge.
- [5] Chen, Y. (1998) Arithmetic circuit verification based on word-level decision diagrams. PhD Thesis, School of Computer Science, Carnegie Mellon University, Pittsburgh, USA.
- [6] IEEE STD-754 (1985) *IEEE standard for binary floating-point arithmetic*. The Institute of Electrical and Electronics Engineers, New York, USA.
- [7] Harrison, J.R. (2000) Floating-point verification in HOL light: the exponential function. *Form. Methods Syst. Des.* **16**, 271–305.
- [8] Tang, P.T.P. (1989) Table-driven implementation of the exponential function in IEEE floating-point arithmetic. *ACM Trans. Math. Softw.*, **15**, 144–157.
- [9] Bui, H.T., Khalaf, B. and Tahar, S. (1999) Table-Driven Floating-Point Exponential Function. *Proc. CCECE 99*, Edmonton, AB, Canada, May 9–12, pp. 450–455. IEEE Canada, Dundas, ON, Canada.
- [10] Abdel-Hamid, A.T., Tahar, S. and Harrison, J. (2002) Enabling Hardware Verification Through Design Changes. *Proc. ICFEM 02*, Shanghai, China, October 21–25, *Lecture Notes in Computer Science*, Vol. 2495, pp. 459–470. Springer, Berlin.
- [11] Barrett, G. (1989) Formal methods applied to a floating point number system. *IEEE Trans. Softw. Eng.*, **SE-15**, 611–621.
- [12] Miner, P.S. (1995) Defining the IEEE-854 floating-point standard in PVS. Technical Memorandum 110167, NASA Langley Research Center, Hampton, VA 23681-0001, USA.
- [13] IEEE STD-854 (1987) *IEEE standard for radix-independent floating-point arithmetic*. The Institute of Electrical and Electronics Engineers, New York, USA.
- [14] Miner, P.S. and Leathrum, J.F. (1996) Verification of IEEE Compliant Subtractive Division Algorithms. *Proc. FMCAD 96*, Palo Alto, CA, November 6–8, *Lecture Notes in Computer Science*, Vol. 1166, pp. 64–78. Springer, Berlin.
- [15] Carreño, V.A. (1995) Interpretation of IEEE-854 Floating-Point Standard and Definition in the HOL system, Technical Memorandum 110189, NASA Langley Research Center, Hampton, VA 23681-0001, USA.
- [16] Harrison, J.R. (1994) Constructing the real numbers in HOL. *Form. Methods Syst. Des.* **5**, 35–59.
- [17] Harrison, J.R. (1999) A Machine-Checked Theory of Floating-Point Arithmetic. *Proc. TPHOLs 99*, Nice, France, September 14–17. *Lecture Notes in Computer Science*, Vol. 1690, pp. 113–130. Springer, Berlin.
- [18] Harrison, J.R. (2000) Formal Verification of Floating Point Trigonometric Functions. *Proc. FMCAD 00*, Austin, TX, USA, November 1–3. *Lecture Notes in Computer Science*, Vol. 1954, pp. 217–233. Springer, Berlin.
- [19] Harrison, J.R. (2000) Formal Verification of IA-64 Division Algorithms. *Proc. TPHOLs 00*, Portland, OR, August 14–18. *Lecture Notes in Computer Science*, Vol. 1869, pp. 234–251. Springer, Berlin.
- [20] Moore, J.S., Lynch, T. and Kaufmann, M. (1998) A mechanically checked proof of the correctness of the kernel of the AMD5K86 floating-point division algorithm. *IEEE Trans. Comput.* **47**, 913–926.
- [21] Russinoff, D.M. (2000) A Case Study in Formal Verification of Register-Transfer Logic With ACL2: The Floating-Point Adder of the AMD Athlon Processor. *Proc. FMCAD 00*, Austin, TX, USA, November 1–3. *Lecture Notes in Computer Science*, Vol. 1954, pp. 3–36. Springer, Berlin.
- [22] Leaser, M. and O’Leary, J. (1995) Verification of a Subtractive Radix-2 Square Root Algorithm and Implementation. *Proc. ICCD 95*, Austin, TX, USA, October 2–4, pp. 526–531. IEEE Computer Society, Washington, DC, USA.
- [23] Aagaard, M.D. and Seger, C.-J.H. (1995) The Formal Verification of a Pipelined Double-Precision IEEE Floating-Point Multiplier. *Proc. ICCD 95*, San Jose, CA, USA, November 5–9, pp. 7–10. IEEE Computer Society, Washington, DC, USA.
- [24] O’Leary, J., Zhao, X., Gerth, R. and Seger, C.-J.H. (1999) Formally verifying IEEE compliance of floating-point hardware. *Intel Technol. J.*, **Q1**, 1–14.
- [25] Chen, Y.A. and Bryant, R.E. (1998) Verification of Floating Point Addresss. *Proc. CAV 98*, Vancouver, BC, 28 June–2 July, *Lecture Notes in Computer Science*, Vol. 1427, pp. 488–499. Springer, Berlin.



- [26] Cornea-Hasegan, M. (1998) Proving the IEEE correctness of iterative floating-point square root, divide, and remainder algorithms. *Intel Technol. J.*, **Q2**, 1–11.
- [27] Daumas, M., Rideau, L. and Théry, L. (2001) A Generic Library for Floating-Point Numbers and its Application to Exact Computing. *Proc. TPHOLs 01*, Edinburgh, Scotland, September 3–6, *Lecture Notes in Computer Science*, Vol. 2152, pp. 169–184. Springer, Berlin.
- [28] Boldo, S., Daumas, M. and Théry, L. (2003) Formal Proofs and Computations in Finite Precision Arithmetic. *Proc. CALCULEMUS 03*, Rome, Italy, September 10–12, pp. 101–111.
- [29] Boldo, S. and Daumas, M. (2004) Properties of two's complement floating point notations. *Int. J. Softw. Tools Technol. Transf.* **5**, 237–246.
- [30] Berg, C. and Jacobi, C. (2001) Formal Verification of the VAMP Floating Point Unit. *Proc. CHARME 01*, Livingston, Scotland, September 4–7, *Lecture Notes in Computer Science*, Vol. 2144, pp. 325–339. Springer, Berlin.
- [31] Mueller, S.M. and Paul, W.J. (2000) *Computer Architecture. Complexity and Correctness*. Springer, Berlin.
- [32] Beyer, S., Jacobi, C., Kröning, D., Leinenbach, D. and Paul, W.J. (2003) Instantiating Uninterpreted Functional Units and Memory System: Functional Verification of the VAMP. *Proc. CHARME 03*, L'Aquila, Italy, October 21–24, *Lecture Notes in Computer Science*, Vol. 2860, pp. 51–65. Springer, Berlin.
- [33] Sawada, J. and Gamboa, R. (2002) Mechanical Verification of a Square Root Algorithm Using Taylor's Theorem. *Proc. FMCAD 02*, Portland, OR, November 6–8, *Lecture Notes in Computer Science*, Vol. 2517, pp. 274–291. Springer, Berlin.
- [34] Kaivola, R. and Aagaard, M.D. (2000) Divider Circuit Verification With Model Checking and Theorem Proving. *Proc. TPHOLs 00*, Portland, OR, August 14–18, *Lecture Notes in Computer Science*, Vol. 1869, pp. 338–355. Springer, Berlin.
- [35] Kaivola, R. and Kohatsu, K.R. (2003) Proof engineering in the large: formal verification of Pentium® 4 floating-point divider. *Int. J. Softw. Tools Technol. Transf.* **4**, 323–334.
- [36] Kaivola, R. and Narasimhan, N. (2002) Formal Verification of the Pentium® 4 Floating-Point Multiplier. *Proc. DATE 02*, Paris, France, March 4–8, pp. 20–27.
- [37] Melham, T. (1993) *Higher-Order Logic and Hardware Verification*. Cambridge University Press, Cambridge.
- [38] Baier, C. and Katoen, J.-P. (2008) *Principles of Model Checking*. The MIT Press, Cambridge.
- [39] Li, R.-C. (2004) Near optimality of Chebyshev interpolation for elementary function computations. *IEEE Trans. Comput.* **53**, 678–687.
- [40] Brisebarre, N., Muller, J.-M. and Tisserand, A. (2006) Computing machine-efficient polynomial approximations. *ACM Trans. Math. Softw.* **32**, 236–256.

## APPENDIX A. DETAILS OF RTL TO ALGORITHMIC VERIFICATION

*Verification of addition block.* Figure A1 shows the block diagram of the addition function. The addition procedure covers both the addition and the subtraction operations. The idea is mainly the same for both but handling both cases together brings an added degree of complexity. The algorithm puts both numbers to the same exponent, adds or subtracts the numbers and then normalizes. The first part of the addition procedure checks which input is greater (*onebigger*). This is especially important in cases where the inputs are of opposite signs. If the inputs carry the same sign, the output sign will then be the same. When the signs are different, the input with the greater magnitude will impose its sign. The next step is to denormalize both inputs and perform the addition. However, before going on to that step, '01' has to be concatenated to both numbers (*mbuff1*, *mbuff3*). The reason for this is that the 1 is the *implicit* 1 contained in the IEEE-754 format. The 0 is there to make sure that the carry bit is not lost. Denormalizing is done by right-shifting the smaller input by an amount determined by the difference in exponents (*Counter*). The exponent is unbiased by removing 127 ('0111111') from its biased value (*mbuff6*). Addition is then performed normally and the last part is normalizing. It would have been more convenient to use FOR loops for denormalizing purposes but the code would have been more dense and significantly more complex.

In HOL, we modeled this algorithm as follows:

```

 $\vdash_{def}$  ADDER1_RTL s1 s2 s3 e1 e2 e3 m1 m2 m3 =
   $\exists$  onebigger counter count mbuff1 mbuff2 mbuff3 mbuff4 mbuff5 mbuff6.
    (if BINVAL e1 > BINVAL e2 then onebigger = T
     else
       (if BINVAL e2 > BINVAL e1 then onebigger = F
        else
          (if BINVAL m1 > BINVAL m2 then onebigger = T
           else onebigger = F)))  $\wedge$ 
    (if s1 = s2 then
      s3 = s1
    else
      (if onebigger = T then s3 = s1 else s3 = s2))  $\wedge$ 
    (if onebigger = F then
      (counter = BINVAL e2 - BINVAL e1)  $\wedge$ 
      (mbuff1 = WCAT (WORD [F; T], m1))  $\wedge$ 
      (mbuff3 = WCAT (WORD [F; T], m2))  $\wedge$ 
      (mbuff6 = BINVAL e2 - 127))

```



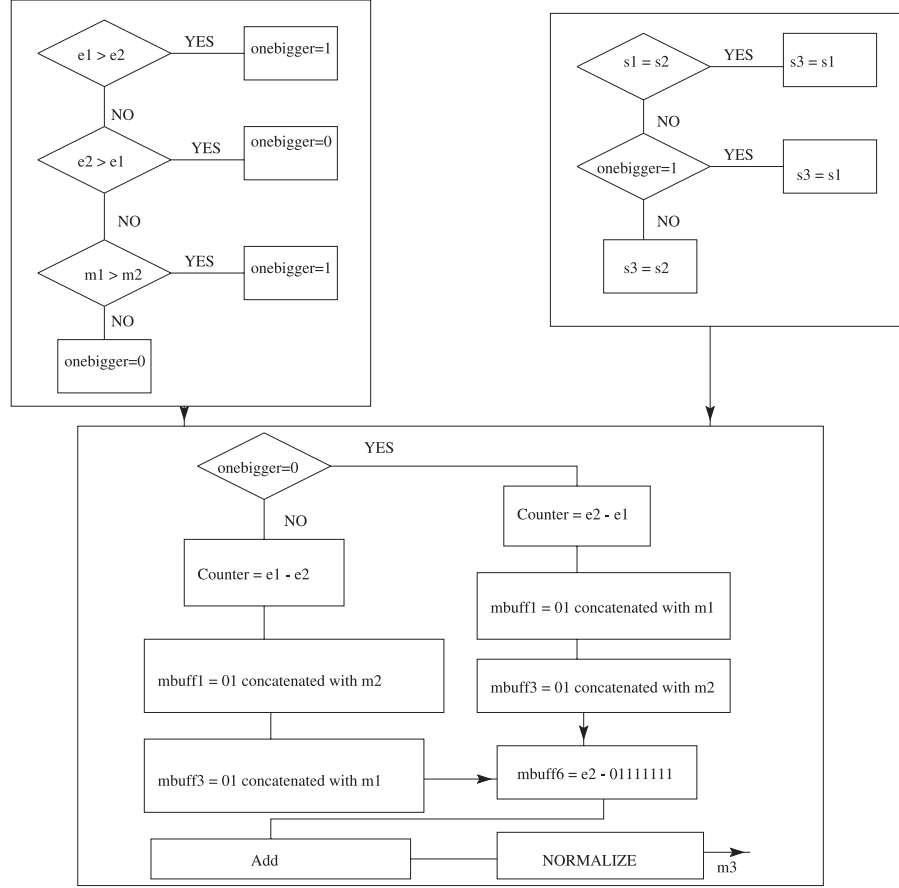


FIGURE A1. Addition block diagram.

```

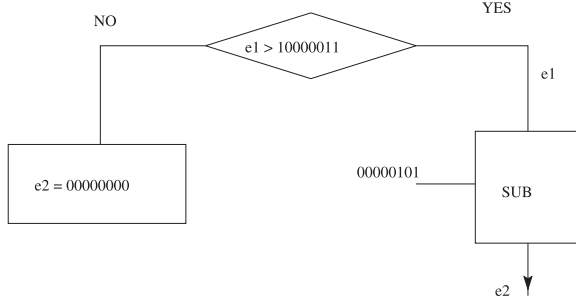
else
  (counter = BNVAL e1 - BNVAL e2) ^
  (mbuff1 = WCAT (WORD [F; T], m2)) ^
  (mbuff3 = WCAT (WORD [F; T], m1)) ^
  (mbuff6 = BNVAL e1 - 127) ^
  (mbuff2 = SHRN mbuff1 counter) ^
  (if s1 = s2 then
    BNVAL mbuff4 = BNVAL mbuff2 + BNVAL mbuff3
  else
    BNVAL mbuff4 = BNVAL mbuff3 - BNVAL mbuff2) ^

  (if BIT 24 mbuff4 = T then
    (mbuff5 = SHLN mbuff4 1) ^ (BNVAL e3 = mbuff6 + 128) ^
    ... ^
  (if BIT 0 mbuff4 = T then
    (mbuff5 = SHLN mbuff4 25) ^ (BNVAL e3 = mbuff6 + 104) ^
  else
    (mbuff5 = WORD (REPLICATE 25 F)) ^ (BNVAL e3 = 0) ^
    (s3 = F) ^ (m3 = WSEG 23 2 mbuff5))

```

where SHRN and SHLN are functions that shifted right and left a Boolean word to a given number of bits, respectively. Following a similar approach to the floating-point multiplier, we proved the following lemma that checks the correctness of the RTL implementation of the floating-point addition function against its algorithmic specification in HOL.

Lemma 4:  $\text{ADDER1\_RTL\_TO\_ALGORITHM\_Correct}$   
 $\vdash \text{ADDER1\_RTL } s1 \ s2 \ s3 \ e1 \ e2 \ e3 \ m1 \ m2 \ m3 \Rightarrow$   
 $(\text{float } (BV \ s3, BNVAL \ e3, BNVAL \ m3)) =$   
 $\text{float } (BV \ s1, BNVAL \ e1, BNVAL \ m1) +$   
 $\text{float } (BV \ s2, BNVAL \ e2, BNVAL \ m2))$



**FIGURE A2.** Division by 32 block diagram.

*Verification of division by 32 block.* This function is only required to be used on a specific type of numbers: multiples of 32. Knowing this fact, the procedure does not need to support all possible ranges of inputs. The operations performed can be explained as follows: the algorithm will output zero if the input exponent is  $< 5$  and will simply subtract five from the exponent if it is not the case. This can be seen in Fig. A2.

In HOL, we modeled this algorithm as follows:

```

 $\vdash_{def}$  D32_RTL s1 s2 e1 e2 m1 m2 =
  (s2 = s1)  $\wedge$ 
  (if BINVAL e1 > 131 then
    (e2 = NBWORD 8 (BINVAL e1 - 5))  $\wedge$ 
    (m2 = m1)
  else
    (e2 = NBWORD 8 0)  $\wedge$  (m2 = NBWORD 23 0))
  
```

Then we established the correctness of the algorithmic to RTL transition of the division by 32 block in HOL as follows:

```

 $\vdash_{def}$  ROUND1_RTL s1 s2 e1 e2 m1 m2 =
   $\exists$  mbuff1 imp1 imp2 imp3 imp4 imp5 imp6 imp7 count.
  (if BINVAL e1 < 126 then
    (e2 = NBWORD 8 0)  $\wedge$  (m2 = NBWORD 23 0)  $\wedge$  (s2 = F)
  else
    (if BINVAL e1 = 126 then
      (e2 = NBWORD 8 127)  $\wedge$  (m2 = NBWORD 23 0)  $\wedge$  (s2 = s1)
    else
      (s2 = s1)  $\wedge$  (mbuff1 = BINVAL e1 - 127)  $\wedge$ 
      (imp1 = WCAT (WORD [F], m1))  $\wedge$ 
      (count = 23 - mbuff1)  $\wedge$ 
      (imp2 = SHRN imp1 (count - 1))  $\wedge$ 
      (imp3 = BIT 0 imp2)  $\wedge$ 
      (imp4 = SHRN imp2 1)  $\wedge$ 
      (imp5 = NBWORD 24 (BINVAL imp4 + BV imp3))  $\wedge$ 
      (if mbuff1 < 23 then
        (if BIT (mbuff1 + 1) imp5 = T then
          (BINVAL e2 = mbuff1 + 1 + 127)  $\wedge$ 
          (imp6 = SHLN imp5 (count - 1))
        else
          (BINVAL e2 = mbuff1 + 127)  $\wedge$ 
          (imp6 = SHLN imp5 count))  $\wedge$ 
          (m2 = WSEG 23 0 imp6)
        else
          (e2 = e1)  $\wedge$  (m2 = m1))))
  
```

Lemma 5: D32\_RTL\_TO\_ALGORITHM\_Correct

```

 $\vdash$  D32_RTL s1 s2 e1 e2 m1 m2  $\implies$ 
  (Toint (float (BV s2, BINVAL e2, BINVAL m2)) =
   Toint (float (BV s1, BINVAL e1, BINVAL m1)) / Int_32)
  
```

where Toint is the coercion for mapping floating-point numbers into machine integers, and ‘/’ denotes the division operation on 2’s complement 32 bit machine integers. The lemma is proved by rewriting on definitions of the division function in RTL (D32) and algorithmic (‘/’) level, and valuations on floating-point numbers and machine integers.

*Verification of round to nearest integer block.* Figure A3 illustrates the round to nearest algorithm. It starts by checking if the exponent is of the order of  $-2$  or less. This would result in an output of zero. The second case is to check if the exponent is  $-1$  in which case the output would be equal to 1. These two IF statements are for negative exponent handling since the main algorithm cannot deal with these cases.

The basic idea is to verify the bit at the 0.5 position. If the bit is set, the decimal positions are filled with zero and we add one to the resulting integer. If the bit is reset, the bits located to the right of the decimal point will be reset. To accomplish this, the input is first shifted right by a number of positions corresponding to the exponent (so that all fraction bits are shifted out). The number obtained should be an integer. This number is then incremented by one if the bit at 0.5 is set else it should be left the same.

In HOL we modeled this algorithm as follows:

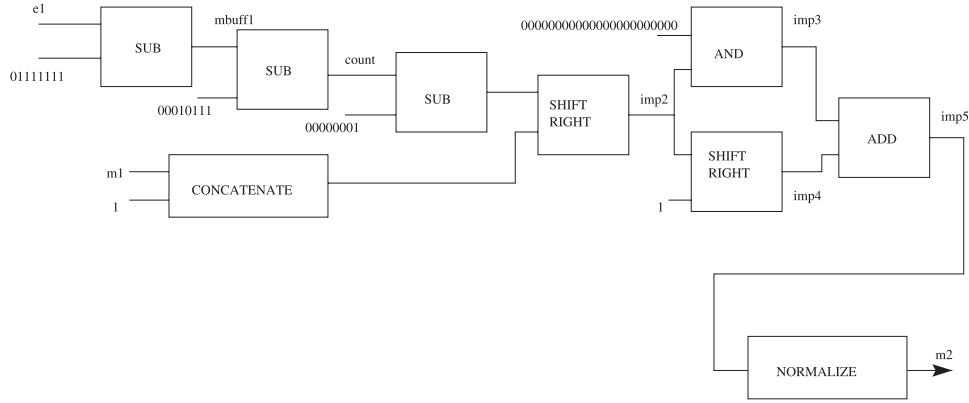


FIGURE A3. Round to nearest integer block diagram.

Then we established the correctness of the algorithmic to RTL transition of the round to nearest integer block in HOL as follows:

Lemma 6: ROUND1\_RTL\_TO\_ALGORITHM\_Correct  
 $\vdash \text{ROUND1\_RTL } s1 \ s2 \ e1 \ e2 \ m1 \ m2 \implies$   
 $(\text{Toint } (\text{float } (\text{BV } s2, \text{BNVAL } e2, \text{BNVAL } m2))) =$   
 $\text{INTRND } (\text{float } (\text{BV } s1, \text{BNVAL } e1, \text{BNVAL } m1))$

where the function INTRND is the composition of the round-to-integer-value operation on floating-point numbers (ROUND FLOAT) and the coercion function Toint. Therefore, our main task in the proof of the above-mentioned theorem was to show that the result of the rounding operation following the RTL algorithm is the best approximation to the real value of the input number. This is established in HOL as the following theorem:

Lemma 7:  
 $\vdash \text{ROUND1\_RTL } s1 \ s2 \ e1 \ e2 \ m1 \ m2 \implies$   
 $((\text{BV } s2, \text{BNVAL } e2, \text{BNVAL } m2) =$   
 $\text{closest } (\text{valof float\_format}) \ (\lambda a. \text{T})$   
 $\{a \mid \text{is\_integral float\_format } a\}$   
 $(\text{valof float\_format } (\text{BV } s1, \text{BNVAL } e1, \text{BNVAL } m1)))$

$\vdash_{\text{def}} \text{MOD32\_RTL } s1 \ s2 \ e1 \ e2 \ m1 \ m2 =$   
 $\exists \text{mbuff1 imp1 imp2 imp3 imp4 imp5 imp6 count n2 n3.}$   
 $(s2 = s1) \wedge (\text{BNVAL mbuff1} = \text{BNVAL } e1 - 127) \wedge$   
 $(\text{imp1} = \text{WCAT } (\text{WORD } [T], m1)) \wedge$   
 $(\text{if BNVAL mbuff1} > 23 \text{ then}$   
 $\quad (\text{count} = \text{BNVAL mbuff1} - 23) \wedge (\text{imp2} = \text{SHLN imp1 count})$   
 $\text{else}$   
 $\quad (\text{count} = 23 - \text{BNVAL mbuff1}) \wedge (\text{imp2} = \text{SHRN imp1 count})) \wedge$   
 $(\text{if BIT 7 mbuff1} = T \text{ then}$   
 $\quad n2 = \text{NBWORD } 24 \ 0$   
 $\text{else}$   
 $\quad n2 = \text{NBWORD } 24 \ 31 \ \text{WAND } \text{imp2}) \wedge$   
 $(\text{if BIT 4 n2} = T \text{ then}$   
 $\quad (\text{BNVAL } e2 = 131) \wedge (n3 = \text{SHLN } n2 \ 19)$   
 $\text{else}$   
 $\quad (\text{if BIT 3 n2} = T \text{ then}$   
 $\quad \quad (\text{BNVAL } e2 = 130) \wedge (n3 = \text{SHLN } n2 \ 20)$

where is\_integral checks if a floating-point number has a finite and integer value.

*Verification of Modulo 32 block.* Modulo 32 is an operation that is done by simply taking the five first bits located to the left of the decimal point. The result will then be an unsigned 5-bit integer that will have to be converted to the single-precision format. For the block diagram, refer to Figure A4. The lower right portion of the figure shows a loop that was not actually implemented in the algorithm. It was drawn like that in order to reduce the complexity of the diagram. The variable 'I' is used as a loop variable.

The procedure is somewhat similar to that of rounding to the nearest integer. The input is first shifted right by the number of bits corresponding to the exponent. The result is then ANDed with the '1111' bit pattern in order to isolate the five bits. The conversion process checks where the first 1 is located starting from the most significant position. An exponent is then assigned accordingly and the result is shifted left to comply with the rules of normalization.

In HOL we modeled this algorithm as follows:



```

else
  (expo = WORD [F; F]) ∧
  (if BNVAL m1 > BNVAL m2 then
    magn = WORD [T; F]
  else
    (if BNVAL m2 > BNVAL m1 then
      magn = WORD [F; T]
    else
      magn = WORD [F; F])) ∧
  (if sign = WORD [F; F] then
    (if expo = WORD [T; F] then
      flag = WORD [T; F; F]
    else
      (if expo = WORD [F; T] then
        flag = WORD [F; F; T]
      else
        (if magn = WORD [T; F] then
          flag = WORD [T; F; F]
        else
          (if magn = WORD [F; T] then
            flag = WORD [F; F; T]
          else
            flag = WORD [F; T; F]))))
else
  (if sign = WORD [T; T] then
    (if expo = WORD [T; F] then
      flag = WORD [F; F; T]

```

```

else
  (if expo = WORD [F; T] then
    flag = WORD [T; F; F]
  else
    (if magn = WORD [T; F] then
      flag = WORD [F; F; T]
    else
      (if magn = WORD [F; T] then
        flag = WORD [T; F; F]
      else
        flag = WORD [F; T; F]))))
else
  (if sign = WORD [T; F] then
    flag = WORD [F; F; T]
  else
    flag = WORD [T; F; F]))

```

Figure A5 shows the block diagram of the comparison function.

The correctness of the algorithmic to RTL transition of the comparison block is established in HOL as follows:

Lemma 9: COMP\_RTL\_TO\_ALGORITHM\_Correct

```

⊢ COMP_RTL s1 s2 e1 e2 m1 m2 flag ⇒
  (if flag = WORD [T; F; F] then
    Toint (float (BV s1, BNVAL e1, BNVAL m1)) >
    Toint (float (BV s2, BNVAL e2, BNVAL m2))
  else
    (if flag = WORD [F; T; F] then
      Toint (float (BV s1, BNVAL e1, BNVAL m1)) =

```

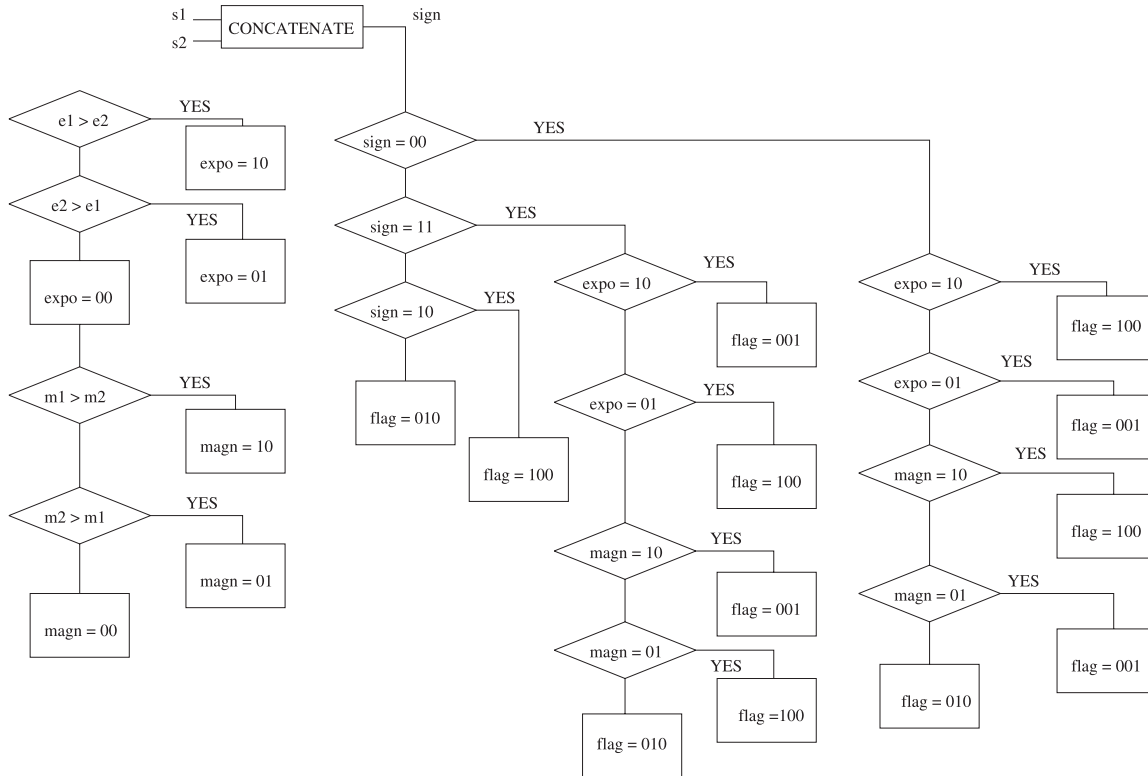


FIGURE A5. Comparison block diagram.

```
Toint (float (BV s2,BNVAL e2,BNVAL m2))
  else
    Toint (float (BV s1,BNVAL e1,BNVAL m1)) <
      Toint (float (BV s2,BNVAL e2,BNVAL m2)))
```

Note that we use conventional symbols for comparison operations on machine integers in the algorithmic level using the operator overloading feature of HOL. The proof is straightforward by rewriting on the definitions of the functions COMP, and <, =, >, and valuation on machine integers.

*Verification of Powers of Two block.* The powers of two function can be implemented by realizing that the value of the input is the value of the output exponent. For example, placing four as an input would result in two to the power of four, yielding four in the exponent field. The objective of the function would then be to convert the input, being an IEEE-754 number, to a 2's complement number. The bias of 127 would then be added to the result and the sum would be placed in the exponent field. The sign and significand fields will be filled with zeros because the result will always be positive and will always be an integer multiple of two. A detailed block diagram of the Power of Two function is given in Fig. A6. In this figure, there is a loop in the lower left section that is used to provide a concise description of the algorithm. It uses the variable 'I' as a loop variable.

In HOL we modeled this algorithm as follows:

```
⊢def TWOPOWERM_RTL s1 e1 m1 s3 e3 m3 =
  ∃ expo magn buff buff2.
    (if e1 = WORD [F; F; F; F; F; F; F; F] then
      (e3 = WORD [F; T; T; T; T; T; T; T]) ∧
      (m3 = WORD (REPLICATE 23 F)) ∧ (s3 = F)
    else
      ((expo = NBWORD 8 (BNVAL e1 - 127)) ∧
       (magn = WCAT (WORD [F; F; F; F; F; F; F; T],m1))) ∧
      (if s1 = F then
        (if expo = WORD [F; F; F; F; F; F; T; T] then
          buff = WSEG 8 16 magn
        else
          (if expo = WORD [F; F; F; F; F; F; T; F] then
            buff = WSEG 8 17 magn
          else
            (if ... else
              buff = WSEG 8 23 magn))))))
    else
      (if expo = WORD [F; F; F; F; F; F; T; T] then
        buff2 = WSEG 8 16 magn
      else
        (if expo = WORD [F; F; F; F; F; F; T; F] then
          buff2 = WSEG 8 17 magn
        else
          (if ... else
            buff2 = WSEG 8 23 magn))))))
      (buff = NBWORD 8 (0 - BNVAL buff2))) ∧ (s3 = F) ∧
      (e3 = NBWORD 8 (BNVAL buff + 127)) ∧ (m3 = NBWORD 23 0))
```

```
Lemma 10: TWOPOWERM_RTL_TO_ALGORITHM_Correct
⊢ TWOPOWERM_RTL s1 e1 m1 s3 e3 m3 ⇒
  (valof float_format (BV s3,BNVAL e3,BNVAL m3) =
   exp (Ival (Toint (float
    (BV s1,BNVAL e1,BNVAL m1))) * ln 2))
```

where Ival is the valuation function on machine integers, and exp and ln are the exponential and natural logarithmic functions defined in the HOL real library, respectively. The proof is done by rewriting on the definition of the function TWOPOWERM and valuations on floating-point numbers and machine integers.

*Verification of Get J block.* The current implementation of the exponential circuit is the table-driven implementation. The table index should ideally be an unsigned integer to make the search easier. The 'Get J' procedure takes care of this. It takes a number in the single-precision format and transforms it to an unsigned number. The procedure examines the exponent and extracts the corresponding bits from the significand. Even though the source code uses a series of IF statements, the block diagram in Fig. A7 shows a loop that uses a variable 'I' to perform the required task.

Using an unsigned number for the search makes the task of finding a correct value for S easier (refer to the algorithm described in Section 3).

The correctness of the algorithmic to RTL transition of the powers of two block is established in HOL as follows:



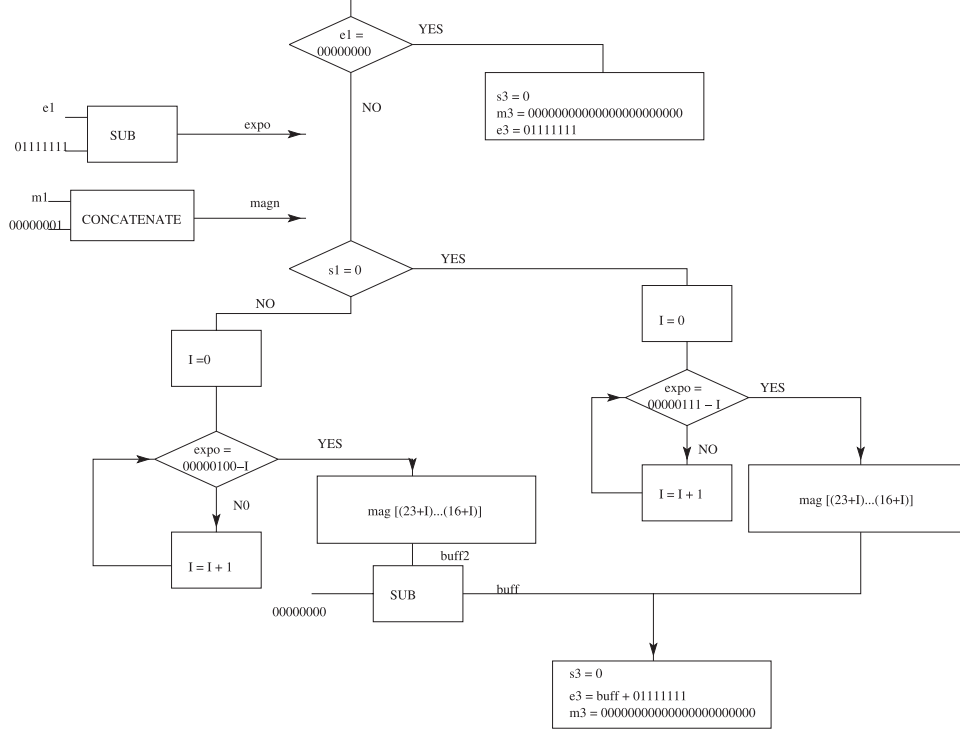


FIGURE A6. Powers of Two block diagram.

In HOL we modeled this algorithm as follows:

```

 $\vdash_{def}$  GET_J_RTL s1 e1 m1 j =
   $\exists$  expo magn.
  (BNVAL expo = BNVAL e1 - 127)  $\wedge$ 
  (magn = WCAT (WORD [F;F;F;F;T], m1))  $\wedge$ 
  (if expo = WORD [F;F;F;F;F;T;F;F] then
    j = WSEG 5 19 magn
  else
    (if expo = WORD [F;F;F;F;F;F;T;T] then
      j = WSEG 5 20 magn
    else
      (if expo = WORD [F;F;F;F;F;F;T;F] then
        j = WSEG 5 21 magn
      else
        (if expo = WORD [F;F;F;F;F;F;F;T] then
          j = WSEG 5 22 magn
        else
          (if expo = WORD [F;F;F;F;F;F;F;F] then
            j = WSEG 5 23 magn
          else
            j = WORD (REPLICATE 5 F))))))
    
```

The correctness of the algorithmic to RTL transition of the get J block is established in HOL as follows:

```

Lemma 11: GET_J_RTL_TO_ALGORITHM_Correct
 $\vdash$  GET_J s1 e1 m1 j  $\implies$  (Int (BNVAL j) =
  Toint (float (BV s1, BNVAL e1, BNVAL m1)))
    
```

where Int is the bijection function that converts a natural number to the machine integer type. The proof is done by

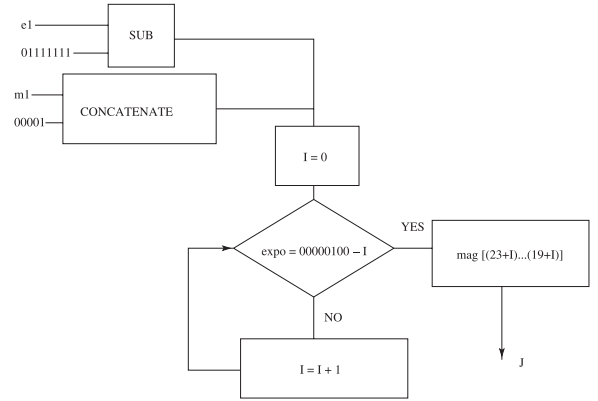


FIGURE A7. Get J block diagram.

rewriting on the definition of the function GET\_J, case analysis on the input exponent and valuations on floating-point numbers and machine integers.

## APPENDIX B. DETAILS OF GATE LEVEL TO RTL VERIFICATION

*Verification of n-bit Adder.* The n-bit Adder in the RTL level is specified as follows:

```

 $\vdash_{def}$  nadd_spec n a b cin s c =
  (BNVAL s = ((if ((BNVAL a + BNVAL b + BV cin) < (2 EXP (SUC n))) then
    (BNVAL a + BNVAL b + BV cin)
  else
    ((BNVAL a + BNVAL b + BV cin) - 2 EXP (SUC n)))))  $\wedge$ 
  (c =  $\neg$ ((BNVAL a + BNVAL b + BV cin) < (2 EXP (SUC n))))

```

The  $n$ -bit Adder at the gate level is implemented using primitive building blocks such as AND, OR, NOT and XOR as follows:

```

 $\vdash_{def}$  nadd_imp 0 a b cin sum1 cout =
  fa_imp (a 0) (b 0) cin (sum1 0) cout)  $\wedge$ 
  (nadd_imp (SUC n) a b cin sum1 cout =  $\exists$  (cripple:bool).
  (fa_imp (a (SUC n)) (b (SUC n)) cripple (sum1 (SUC n)) cout)  $\wedge$ 
  (nadd_imp n a b cin sum1 cripple))

```

The correctness of the  $n$ -bit Adder block is proved in HOL as in the following theorem:

Lemma 20: N\_ADD\_GATE\_LEVEL\_TO\_RTL\_Correct  
 $\vdash$  nadd\_imp n a b cin sum cout = nadd\_spec n a b cin sum cout

*Verification of  $n$ -bit Subtractor.* The  $n$ -bit Subtractor in the RTL is specified as follows:

```

 $\vdash_{def}$  nSub_spec 0 a b bin dif bout =
  fs_spec (a 0) (b 0) bin (dif 0) bout)  $\wedge$ 
  (nSub_spec (SUC n) a b bin dif bout =  $\exists$  (briipple:bool).
  (fs_spec (a (SUC n)) (b (SUC n)) briipple (dif (SUC n)) bout)  $\wedge$ 
  (nSub_spec n a b bin dif briipple))

```

The  $n$ -bit Subtractor at the gate level is implemented as follows:

```

 $\vdash_{def}$  (nSub_imp 0 a b bin dif bout =
  fs_imp (a 0) (b 0) bin (dif 0) bout)  $\wedge$ 
  (nSub_imp (SUC n) a b bin dif bout =  $\exists$  (briipple:bool).
  (fs_imp (a (SUC n)) (b (SUC n)) briipple (dif (SUC n)) bout)  $\wedge$ 
  (nSub_imp n a b bin dif briipple))

```

The correctness of the  $n$ -bit Subtractor block is proved in HOL as in the following theorem:

Lemma 21: N\_SUB\_GATE\_LEVEL\_TO\_RTL\_Correct  
 $\vdash$  nSub\_imp n a b bin dif bout  $\implies$  nSub\_spec n a b bin dif bout

*Verification of  $n$ -bit Comparator.* The  $n$ -bit Comparator in the RTL is specified as follows:

```

 $\vdash_{def}$  n_BIT_COMP_Spec n A B lf gf ef l g e =
  (l = ((BNVAL A < BNVAL B)  $\wedge$  ef)  $\vee$  lf)  $\wedge$ 
  (g = ((BNVAL A > BNVAL B)  $\wedge$  ef)  $\vee$  gf)  $\wedge$ 
  (e = ((BNVAL A = BNVAL B)  $\wedge$  ef))

```

The n-bit Comparator at the gate level is implemented as follows:

```

 $\vdash_{def}$  BIT_COMPARE_Imp a b l g e =
   $\exists$  anot bnot s1 s2.
    (not1 a anot)  $\wedge$ 
    (not1 b bnot)  $\wedge$ 
    (and2 anot b l)  $\wedge$ 
    (and2 bnot a g)  $\wedge$ 
    (and2 a b s1)  $\wedge$ 
    (and2 anot bnot s2)  $\wedge$ 
    (or2 s1 s2 e)

 $\vdash_{def}$  n_BIT_COMP_BULID_Imp a b lf gf ef l g e =
   $\exists$  l1 e1 g1 s1 sg.
    (BIT_COMPARE_Imp a b l1 g1 e1)  $\wedge$ 
    (and2 e1 ef e)  $\wedge$ 
    (and2 l1 ef s1)  $\wedge$ 
    (and2 g1 ef sg)  $\wedge$ 
    (or2 s1 lf l)  $\wedge$ 
    (or2 sg gf g)

 $\vdash_{def}$  (n_BIT_COMP_Imp 0 A B lf gf ef l g e =
  n_BIT_COMP_BULID_Imp (A 0) (B 0) lf gf ef l g e)  $\wedge$ 
  (n_BIT_COMP_Imp (SUC n) A B lf gf ef l g e =  $\exists$  l1 g1 e1.
    (n_BIT_COMP_BULID_Imp (A (SUC n)) (B (SUC n)) lf gf ef l1 g1 e1)  $\wedge$ 
    (n_BIT_COMP_Imp n A B l1 g1 e1 l g e))

```

The correctness of the n-bit Comparator block is proved in HOL as in the following theorem:

Theorem: N\_BIT\_COMP\_GATE\_LEVEL\_TO\_RTL\_Correct  
 $\vdash$  n\_BIT\_COMP\_BULID\_Imp a b lf gf ef l g e = n\_BIT\_COMP\_BULID\_Spec a b lf gf ef l g e

*Verification of n-bit Concatenator.* The n-bit Concatenator in the RTL is specified as follows:

```

 $\vdash_{def}$  CONCATENATE_Spec (n:num) (X:num $\rightarrow$ bool) (Y:num $\rightarrow$ bool) =
  (BNVAL Y = 2 EXP (SUC n) + BNVAL X)

```

The n-bit Concatenator at the gate level is implemented as follows:

```

 $\vdash_{def}$  CONCATENATE_Imp (n:num) (X:num $\rightarrow$ bool) (Y:num $\rightarrow$ bool) =
  (Y (SUC n) = T)  $\wedge$ 
  ( $\forall$  n. Y n = X n)

```

The correctness of the n-bit Concatenator block is proved in HOL as in the following theorem:

Theorem: CONCATENATE\_GATE\_LEVEL\_TO\_RTL\_THM  
 $\vdash$  CONCATENATE\_Imp n X Y  $\implies$  CONCATENATE\_Spec n X Y

*Verification of n-bit Multiplexer.* The n-bit Multiplexer in the RTL is specified as follows:

```

 $\vdash_{def}$  MUX_Spec n A B s OUT =
  (BNVAL OUT = (if (s = F) then BNVAL A else BNVAL B))

```

The n-bit Multiplexer at the gate level is implemented as follows:

```

 $\vdash_{def}$  MUX_Cell_Imp a b s out =
   $\exists$  s1 s2 s3.
    (not1 s s1)  $\wedge$ 
    (and2 s1 a s2)  $\wedge$ 
    (and2 s b s3)  $\wedge$ 
    (or2 s2 s3 out)

 $\vdash_{def}$  (MUX_Imp 0 A (B:num $\rightarrow$  bool) s OUT =
  MUX_Cell_Imp (A 0) (B 0) s (OUT 0))  $\wedge$ 
  (MUX_Imp (SUC n) A B s OUT =
  (MUX_Cell_Imp (A (SUC n)) (B (SUC n)) s (OUT (SUC n)))  $\wedge$ 
  (MUX_Imp n A (B:num $\rightarrow$  bool) s OUT))

```

The correctness of the n-bit Multiplexer block is proved in HOL as in the following theorem:

Theorem: MUX\_GATE\_LEVEL\_TO\_RTL\_Correct  
 $\vdash$  MUX\_Imp n A B s OUT = MUX\_Spec n A B s OUT

*Verification of n-bit Shifter.* The n-bit Shifter in the RTL is specified as follows:

```

 $\vdash_{def}$  ShiftRight_Spec n M L =
  ((BNVAL L * (2 EXP n)) + BNVAL (WSEG n 0 M) = (BNVAL M))

 $\vdash_{def}$  ShiftLeFT_Spec n M L =
  (BNVAL L = (2 EXP n) * BNVAL M)

```

The n-bit Shifter at the gate level is implemented as follows:

```

 $\vdash_{def}$  (ShiftLeFT_Imp 0 M L =
  ((L 1) = (M 0))  $\wedge$ 
  (L 0 = F) )  $\wedge$ 
  (ShiftLeFT_Imp (SUC n) M L =
  (L (SUC (SUC n)) = M (SUC n))  $\wedge$ 
  ShiftLeFT_Imp n M L)

 $\vdash_{def}$  (ShiftRight_Imp 0 M L =
  ((L 0) = (M 1))  $\wedge$ 
  (L 1 = F) )  $\wedge$ 
  (ShiftRight_Imp (SUC n) M L =
  (L (SUC n) = M (SUC (SUC n)))  $\wedge$ 
  ShiftRight_Imp n M L)

```

The correctness of the n-bit Shifter block is proved in HOL as in the following theorems:

Theorem: SHIFT\_LEFT\_GATE\_LEVEL\_TO\_RTL\_THM  
 $\vdash$  ShiftLeFT\_Imp n M L  $\implies$  ShiftLeFT\_Spec n M L

Theorem: SHIFT\_RIGHT\_GATE\_LEVEL\_TO\_RTL\_THM  
 $\vdash$  ShiftRight\_Imp n M L  $\implies$  ShiftRight\_Spec n M L